

THE LEUCOCYTE ANTIGEN *FactsBook*

A. Neil Barclay
Marian L. Birkeland
Marion H. Brown
Albertus D. Beyers
Simon J. Davis
Chamorro Somoza
Alan F. Williams

*MRC Cellular Immunology Unit
Sir William Dunn School of Pathology
University of Oxford, Oxford, UK*



Academic Press
Harcourt Brace Jovanovich, Publishers
LONDON SAN DIEGO NEW YORK BOSTON
SYDNEY TOKYO TORONTO

459079



This book is printed on acid-free paper

ACADEMIC PRESS LIMITED
24-28 Oval Road
LONDON NW1 7DX

United States Edition published by
ACADEMIC PRESS INC.
San Diego, CA 92101

Copyright © 1993 by
ACADEMIC PRESS LIMITED

All rights reserved

No part of this book may be reproduced in any form by photostat, microfilm, or by
any other means, without written permission from the publishers

A catalogue record for this book is available from the British Library

ISBN 0-12-178180-1

Designed and typeset by Eric Drewery and Adrian Singer

Printed and bound in Great Britain by Mackays of Chatham
Chatham, Kent

Preface

Abbreviations

Dedications

Section

Chapter

Introduction

Chapter

The Analysis

Chapter

Protein

Chapter

Chromatography

Section

CD model

CD1

CD2

CD3/4

CD4

CD5

CD6

CD7

CD8

CD9

CD10

CD11

CD12

CD13

CD14

CD15

CD16

CDw1

CD18

CD19

CD20

CD21

CD22

CD23

CD24

CD25

CD26

Contents

Preface	VII
Abbreviations	VIII
Dedication	IX

Section I THE INTRODUCTORY CHAPTERS

<i>Chapter 1</i>	
Introduction	2
<i>Chapter 2</i>	
The Analysis and Architecture of the Leucocyte Cell Surface	13
<i>Chapter 3</i>	
Protein Superfamilies and Cell Surface Molecules	38
<i>Chapter 4</i>	
Chromosomal Organization of the Genes Encoding Leucocyte Surface Antigens	88

Section II THE LEUCOCYTE ANTIGENS

CD molecules

CD1	102	CD27	160	CD53	222
CD2	104	CD28	162	CD54	224
CD3/TcR	106	CD29	164	CD55	226
CD4	110	CD30	166	CD56	228
CD5	112	CD31	168	CD57	231
CD6	114	CDw32	170	CD58	232
CD7	116	CD33	174	CD59	234
CD8	118	CD34	176	CDw60	237
CD9	120	CD35	178	CD61	238
CD10	122	CD36	182	CD62	240
CD11	124	CD37	184	CD63	242
CD12	129	CD38	185	CD64	244
CD13	130	CD39	187	CDw65	247
CD14	132	CD40	188	CD66	248
CD15	135	CD41	190	CD67	250
CD16	136	CD42a,b	192	CD68	252
CDw17	139	CD43	196	CD69	254
CD18	140	CD44	198	CDw70	257
CD19	142	CD45	202	CD71	258
CD20	144	CD46	206	CD72	260
CD21	146	CD47	209	CD73	263
CD22	148	CD48	210	CD74	264
CD23	150	CD49	212	CDw75	266
CD24	152	CDw50	217	CD76	267
CD25	154	CD51	218	CD77	268
CD26	158	CDw52	220	CDw78	269

3 Protein Superfamilies and Cell Surface Molecules

CONCEPTS CONCERNING PROTEIN SUPERFAMILIES

Introduction

The amino acid sequences of most leucocyte surface proteins contain segments of sequence that have similarities to other proteins and it is likely that the similar sequences have been derived by divergent evolution from common precursors. Dayhoff et al.¹ introduced the terms "superfamily" for proteins with sequence similarity of 50% or less and "family" for those with more than 50% identity. In protein superfamilies there is often only 15–25% sequence identity and at this level it can be difficult to be confident that a sequence match indicates an evolutionary relationship, rather than just a chance similarity.

The first superfamily of leucocyte surface proteins to be defined was the immunoglobulin superfamily (IgSF) and this is now the largest with more than 100 different polypeptides on a variety of cell types². The sequence identities between the members of this superfamily are at the 15–25% level but analysis showed that the conserved residues are clustered mainly in regions corresponding to the in-pointing residues of β strands of the Ig-fold. In contrast, the regions corresponding to the loops at the ends of the strands mostly show great sequence diversity. It may be regarded as a rule that in superfamilies of sequences that have derived by divergent evolution the conserved residues will relate to important structural features that are characteristic of the superfamily in question.

Several different superfamilies have been identified within leucocyte surface molecules. This chapter describes the methods for their identification and shows alignments of some sequences to illustrate the key residues that are often conserved in these superfamilies. A brief description of the structure and functions of each domain type is given.

Nomenclature for superfamilies, protein domains, repeats and motifs

There is no agreed nomenclature for most superfamilies and thus in this book we have tried to conform to the most commonly used names and in some cases to introduce abbreviations that might be useful. In general where superfamilies are named after a receptor we use the abbreviation "R". For example, the cytokine receptor superfamily is called the "cytokineR" superfamily. This seems useful in that the name becomes distinctive and distinguishes the superfamily usage from discussion in which a receptor is referred to in other ways. The naming of domains is a problem since one might discuss Ig domains either as domains of immunoglobulins or as domains of the superfamily. For the superfamily usage we include the abbreviation "SF" in cases where there may be ambiguities. For example, IgSF, scavengerRSF, FN type IIISF, CCPSF.

The term "domain" is used where it is likely that a segment of sequence forms a discrete structural unit, i.e. a peptide sequence whose three-dimensional conformation is not determined by other parts of the total protein sequence but is "self-contained". Three criteria are considered. First, proof of a domain structure comes from tertiary structure determination. Domains established at this level include: Ig, complement control protein (CCP), EGF, fibronectin (FN) type III, cytokineR and the C-type lectin. The MHC domain has also been revealed by X-ray

crystallography in MHC should not be referred to as an isolated unit rather than domains. However, we segments as domains and superfamily. The folds for are discussed in the comm

Secondly, a domain str that occur as the sole col sequences) within protein containing a high content

A third criterion for d segments are found in th other exons to form a ne variety of structural dom these last two criteria th domains: FN type II, Link

In other cases it is n structural unit, and for t seen with the NGFR sup are always found togeth pattern of exons does not it appears that a precurs repeat, and that addit duplication and divergen independent units withi units interact to form a "repeat" is used, include

The term "motif" is r expected to form a folde protein secretion and C motifs, albeit of rather l good example of a m cytoplasmic domains of signal transduction com 28.

The domains and repe leucocyte surface molec molecules and for surfa kringle, thrombospondi

Identifying domains and There can be problems This is because the le patterns are usually in the superfamily domai program such as FAST, a superfamily that the picks up all superfamily

crystallography in MHC Class I $\alpha 1$ and $\alpha 2$ domains and it might be argued that should not be referred to as a domain since it is not clear that it will be found as an isolated unit rather than appearing always as a structural pair, as for the $\alpha 1$ and $\alpha 2$ domains. However, we will follow precedent in the field and refer to these segments as domains and to the proteins that show this fold as being in the superfamily. The folds for all these domains are illustrated later in this chapter and are discussed in the commentary on each superfamily.

Secondly, a domain structure can also be argued for any superfamily segment that occurs as the sole component of an extracellular sequence, or as sequences (or sequences) within proteins that is contiguous with hinge-like regions of sequence containing a high content of Ala, Gly, Pro, Ser and Thr residues.

A third criterion for defining a sequence as a domain is that the superfamily segments are found in the genome in single exons that can be readily spliced with other exons to form a new gene with an open reading frame. Proteins containing a variety of structural domains could then arise by recombination. On the basis of these last two criteria the following superfamilies can be referred to as containing domains: FN type II, Link, Ly-6, LDLR and the scavengerR.

In other cases it is not clear that a superfamily segment is an independent structural unit, and for these the term "repeat" is used. A good example of this is seen with the NGFR superfamily (Fig. 25), where a block of three or four repeats is always found together without intervening sequence between the repeats. The pattern of exons does not correlate with NGFR repeats. Thus with this superfamily it appears that a precursor structure evolved by gene duplication of a primary repeat, and that additional members of the superfamily have evolved by duplication and divergence of the larger structure. The repeat segments may form independent units within the structure or alternatively it could be that the repeat units interact to form a larger structural unit. Superfamilies for which the term "repeat" is used, include NGFR and leucine-rich glycoprotein repeats.

The term "motif" is used to describe a smaller sequence pattern than a repeat expected to form a folded structural unit. Thus the patterns of signal sequences for protein secretion and GPI attachment (see Chapter 2) would be considered as motifs, albeit of rather ill-defined character in terms of sequence identities. A very good example of a motif is the conserved sequence pattern found in the cytoplasmic domains of the CD3, MB-1, and B29 antigens and other molecules of signal transduction complexes. Alignments identifying this motif are shown in Fig. 28.

The domains and repeats discussed in this chapter include only those present on leucocyte surface molecules. Additional domains have been described for other cell types and for surface molecules of other cell types. These include Fc type I, kringle, thrombospondin, serine protease and perforin domains³.

Identifying domains and repeats: testing the significance of relationships

There can be problems in identifying domains or repeats in new protein sequences. This is because the level of identities is often low and the conserved sequences are usually in small patches throughout the 40–110 residues that make up the superfamily domains and repeats. A first step is to use a database searching program such as FASTA⁴. In many cases this will pick up some of the members of a superfamily that the new protein sequence matches. However, no search program picks up all superfamily members and it is not uncommon for a relationship to be missed.

entirely missed. A second approach is to look by eye, or with a computer program, for the presence of sequence patterns that are characteristic of the different superfamilies. These are thus noted in the sequence line-ups in this chapter. For example, in the IgSF one would look for Cys residues with the patterns L/I/V-X-C and D-X-G-X-Y-X-C for candidate regions that might occur around a conserved disulphide bond. In relation to these there should be other patches, for example V/L/Y-X-W corresponding to β -strand C. If the various conserved patches fall into place then a possible domain has been identified. The candidate domain can be defined in relation to conserved sequence positions and then tested for statistical significance. For example, in the IgSF, the positions of the conserved residues, or equivalent residues if the domain lacks the typical disulphide bond, are nominated and the domain is defined as beginning and ending 20 residues before and after these positions. This proposed domain is then tested for the statistical significance of sequence similarities against a set of domains that are accepted in the IgSF. For other superfamilies, other conserved residues would be chosen and the domain defined in relation to these. Possible key conserved residues are shown in the diagrams in this chapter and the designated residues are used to identify the domains in the entries for the molecules.

In testing for statistical significance of a superfamily relationship it could be argued that the conserved pattern for a domain should be defined and the extent to which this occurs in the new sequence should be tested. However, it is difficult to define precisely a pattern for use in a statistical analysis since many positions in the conserved pattern are one of a group of alternative amino acids. It is difficult to know how to treat sequence gaps in defining a pattern. For example, in the IgSF there can be very large differences in the length of the domain and this creates problems in defining a pattern that is characteristic of the IgSF to use in statistical analysis.

An alternative method to testing a sequence against a single superfamily sequence pattern is to test it against a set of sequences (e.g. 20 sequences) that are accepted as being members of the superfamily in question. In such an analysis a simple statistical program that compares sequences pairwise for similarity can be used and the ALIGN program ¹ has proved satisfactory for this purpose. In these comparisons no account is taken of superfamily patterns. However, if a set of good scores is obtained against a family of sequences, then the superfamily pattern must be present since this is the only pattern in common amongst the family of sequences against which the new domain is being tested.

In the ALIGN program of Dayhoff ¹, the best alignment between two sequences is computed on the basis of a matrix of scores for all possible identities and amino acid substitutions with a penalty scored each time a gap is introduced to improve the number of good matches. The two sequences are then scrambled, realigned and scored again to give a random best score and this is repeated 100–150 times. From the random scores a mean random score plus standard deviation (SD) is calculated and the test score is expressed in terms of its number of SDs above or below the mean random score. Assuming a normal distribution and no effect due to selection of sequences with particular compositions, values of 3.1, 4.3 and 5.5 SD units indicate the probability of the sequence similarity arising by chance is 10^{-3} , 10^{-5} and 10^{-7} respectively. A score of 3 SD is considered to be a threshold for a value that is of interest in indicating a superfamily relationship between proteins. The most common matrix of scores used with the ALIGN program is the

Table 1. The 250 PAM Matrix of scores for sequence matches analysed by the ALIGN program adapted from ref. 107 with permission

C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W	Cys Sulphydryl

Table 1. The 250 PAM Matrix of scores for sequence matches analysed by the ALIGN program adapted from ref. 107 with permission

	C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W	
Cysteine [C]	12	0	-2	-3	-2	-3	-4	-5	-5	-5	-3	-4	-5	-5	-2	-6	-2	-4	0	-8	Cys Sulphydryl
Serine [S]	0	2	1	1	1	1	1	0	0	-1	-1	0	0	-2	-1	-3	-1	-3	-3	-2	Ser
Threonine [T]	-2	1	3	0	1	0	0	0	0	-1	-1	-1	0	-1	0	-2	0	-3	-3	-5	Thr
Proline [P]	-3	1	0	6	1	-1	-1	-1	0	0	0	0	-1	-2	-2	-3	-1	-5	-5	-6	Pro Small hydrophilic
Alanine [A]	-2	1	1	1	2	1	0	0	0	0	-1	-2	-1	-1	-1	-2	0	-4	-3	-6	Ala
Glycine [G]	-3	1	0	-1	1	5	0	1	0	-1	-2	-3	-2	-3	-3	-4	-1	-5	-5	-7	Gly
Asparagine [N]	-4	1	0	-1	0	0	2	2	1	1	2	0	1	-2	-2	-3	-2	-4	-2	-4	Asn
Aspartic [D]	-5	0	0	-1	0	1	2	4	3	2	1	-1	0	-3	-2	-4	-2	-6	-4	-7	Asp Acid, acid amide
Glutamic [E]	-5	0	0	-1	0	0	1	3	4	2	1	-1	0	-2	-2	-3	-2	-5	-4	-7	Glu hydrophilic
Glutamine [Q]	-5	-1	-1	0	0	-1	1	2	2	4	3	1	1	-1	-2	-2	-2	-5	-4	-5	Gln
Histidine [H]	-3	-1	-1	0	-1	-2	2	1	1	3	6	2	0	-2	-2	-2	-2	-2	0	-3	His
Arginine [R]	-4	0	-1	0	-2	-3	0	-1	-1	1	2	6	3	0	-2	-3	-2	-4	-4	2	Arg Basic
Lysine [K]	-5	0	0	-1	-1	-2	1	0	0	1	0	3	5	0	-2	-3	-2	-5	-4	-3	Lys
Methionine [M]	-5	-2	-1	-2	-1	-3	-2	-3	-2	-1	-2	0	0	6	2	4	2	0	-2	-4	Met
Isoleucine [I]	-2	-1	0	-2	-1	-3	-2	-2	-2	-2	-2	-2	-2	2	5	2	4	1	-1	-5	Ile Small hydrophobic
Leucine [L]	-6	-3	-2	-3	-2	-4	-3	-4	-3	-2	-2	-3	-3	4	2	6	2	2	-1	-2	Leu
Valine [V]	-2	-1	0	-1	0	-1	-2	-2	-2	-2	-2	-2	-2	2	4	2	4	-1	-2	-6	Val
Phenylalanine [F]	-4	-3	-3	-5	-4	-5	-4	-6	-5	-5	-2	-4	-5	0	1	2	-1	9	7	0	Phe
Tyrosine [Y]	0	-3	-3	-5	-3	-5	-2	-4	-4	-4	0	-4	-4	-2	-1	-1	-2	7	10	0	Tyr Aromatic
Tryptophan [W]	-8	-2	-5	-6	-6	-7	-4	-7	-7	-5	-3	2	-3	-4	-5	-2	-6	0	0	17	Trp

The mutation matrix is based on the frequency of evolutionary replacements of one amino acid for another at homologous positions between present-day sequences and inferred ancestral sequences. One PAM unit is the unit of evolution represented by the matrix corresponding to one accepted amino acid substitution per 100 residues. This is discussed in detail in ref. 107.

Table 2. ALIGN Scores for comparisons of IgSF domains

V or V-related	$\beta 2$ -microglobulin	Thy-1	N-CAM	CD48	CD45	C or C-related	$\beta 2$ -microglobulin	Thy-1	N-CAM	CD48	CD45
Ig λ	-0.9	7.4	3.3	3.4	-0.6	Ig λ	5.6	1.4	4.7	0.5	-1.0
Ig κ	1.7	3.7	5.4	3.6	-0.1	Ig κ	6.0	1.3	4.0	0.4	-1.5
Ig heavy	1.1	3.9	3.9	4.0	1.5	Ig CH1	4.0	3.0	4.1	1.6	0.7
TcR β	1.8	3.3	4.6	3.6	-0.2	Ig CH2	2.4	2.9	3.8	0.0	0.0
TcR α	2.1	2.3	4.4	3.4	-1.5	Ig CH3	6.3	3.1	3.7	1.4	0.7
TcR γ	-0.2	1.6	3.9	3.3	1.4	TcR β	4.4	2.3	3.0	0.8	2.2
CD8 α	2.6	4.5	4.7	4.3	-0.3	TcR α	2.1	-0.3	1.7	0.3	-0.7
CD4 (1)	2.4	2.5	5.5	4.3	-0.1	TcR γ	1.9	0.8	3.6	1.0	-0.6
PolyIgR (1)	1.5	5.7	2.7	1.9	0.0	MHC I $\alpha 3$	8.2	2.2	2.9	2.4	1.3
PolyIgR (3)	1.7	5.8	4.3	2.4	0.9	MHC II $\alpha 2$	11.2	3.7	4.9	1.6	1.2
MRC OX-2 (1)	-1.1	5.0	5.3	3.3	-0.6	MHC II $\beta 2$	11.3	2.4	4.3	1.1	1.6
Po protein	0.7	3.5	6.0	2.9	-0.4	CD1 $\alpha 3$	9.1	1.3	5.4	0.9	1.2

The IgSF domains were defined from a position 20 residues before the first Cys to 20 residues after the second Cys of a putative Ig-like disulphide bond or equivalent residues in the CD48 sequence. A sequence from rat CD45 is included as a control. This sequence shows none of the conserved patches of sequence characteristic of the IgSF domains but has two Cys and one Trp residue in approximately the same positions as in IgSF domains. The ALIGN program was run with the 250 PAMS Mutation Matrix, a bias of 6 and a break penalty of 6 and 150 random runs were performed. Details of the sequences are given in ref. 108.

mutation matrix calculated replace each other in hom substitution frequency expected for identities between tryptophan and cysteine, this observed to change more frequently. The mutation matrix are illustrated amino acids are favourable evolution. Those that replace chance are shown in bold.

It is stated above that a chance relationship of 10^{-5} comparisons with 20 members superfamily relationship. H and tests have been carried out program in analysis of the Ig patterns were compared with roughly the correct position IgSF patterns ⁶. The results

Table 3. ALIGN Scores for

	GHR
GHR	1
PLR	10.4
GMP130	4.0
EPOR	3.6
IL3Rd1	2.9
GM-CSFR	4.1
IL6R	4.5
IL2R β	2.2
IL4R	1.8
IL3Rd3	4.1
IL7R	1.3

The regions of the domains shown 250 PAMS Mutation Matrix, a scores in SD units are given. putative cytokineRSF domain in

GHR, human growth hormone GMP130, human membrane receptor precursor; IL3Rd1 and human GM-CSF receptor precursor; IL4 precursor.

mutation matrix calculated from the observed frequencies at which amino acids replace each other in homologous proteins between species, compared with the substitution frequency expected by chance. With this matrix a much higher score occurs for identities between residues that rarely change in evolution, such as tryptophan and cysteine, than for small residues like serine and threonine that are observed to change more frequently. The single letter code for amino acids and the mutation matrix are illustrated in Table 1. From this matrix it can be seen which amino acids are favourable substitutes for one another in related proteins in evolution. Those that replace each other more often than would be expected by chance are shown in bold.

It is stated above that an ALIGN score of 4.3 SD indicates a probability of a chance relationship of 10^{-5} and thus it might be thought that one score of 4.3 SD in comparisons with 20 members of a superfamily would alone argue for a reliable superfamily relationship. However, theoretical probabilities may not hold and tests have been carried out to evaluate experimentally the use of the ALIGN program in analysis of the IgSF. Sequences that were thought to have IgSF patterns were compared with others that had two Cys residues and a Trp residue in roughly the correct position for IgSF domains but otherwise did not show typical IgSF patterns ⁶. The results are shown in Table 2 where it can be seen that the

Table 3. ALIGN Scores for alignments of cytokineRSF domains

	GHR	PLR	GMP130	EPOR	IL3Rd1	GM-CSFR	IL6R	IL2Rβ	IL4R	IL3Rd3	IL7R
GHR		10.4	4.0	3.6	2.9	4.1	4.5	2.2	1.8	1.1	1.3
PLR	10.4		8.1	7.2	2.5	7.3	5.2	1.6	1.9	1.4	1.7
GMP130	4.0	8.1		4.7	2.4	4.2	6.2	3.1	2.5	4.5	1.1
EPOR	3.6	7.2	4.7		4.7	5.6	5.1	3.6	2.4	8.7	0.8
IL3Rd1	2.9	2.5	2.4	4.7		4.3	3.1	1.9	3.0	5.3	0.3
GM-CSFR	4.1	7.3	4.2	5.6	4.3		5.3	5.5	2.1	6.7	1.6
IL6R	4.5	5.2	6.2	5.1	3.1	5.3		4.0	1.7	5.2	1.1
IL2Rβ	2.2	1.6	3.1	3.6	1.9	5.5	4.0		2.5	4.3	0.6
IL4R	1.8	1.9	2.5	2.4	3.0	2.1	1.7	2.5		3.8	0.3
IL3Rd3	4.1	4.4	4.5	8.7	5.3	6.7	5.2	4.3	3.8		2.0
IL7R	1.3	1.7	1.1	0.8	0.3	1.6	1.1	0.6	0.3	2.0	

The regions of the domains shown in Fig. 4 were analysed using the ALIGN program with the 250 PAMS Mutation Matrix, a bias of 6, a gap penalty of 6 and 100 random alignments. The scores in SD units are given. Scores of 3 SD or greater are shown in bold. Note for the putative cytokineRSF domain in the IL7 receptor, all scores are less than 3 SD.

GHR, human growth hormone receptor precursor; PLR, rat prolactin receptor precursor; GMP130, human membrane glycoprotein gp130 precursor; EPOR, mouse erythropoietin receptor precursor; IL3Rd1 and d3, mouse IL3 receptor precursor domains 1 and 3; GM-CSFR, human GM-CSF receptor precursor; IL6R, human IL6 receptor precursor; IL2Rβ, human IL2 receptor β chain precursor; IL4R, mouse IL4 receptor precursor; IL7R, human IL7 receptor precursor.

control sequences give an occasional score above 2 SD but that no consistent pattern of good scores is obtained. In contrast the sequences that are considered to belong to the IgSF gave >40% scores of >3 SD. In practice, arguments for an IgSF relationship have proved reliable in terms of subsequent tertiary structure determination in cases where a group of scores >3 SD have been obtained with >33% of one of the sets of IgSF sequences as defined below (i.e. the V set, C1 set or C2 set). Convincing arguments are usually buttressed by one or two scores >5 SD that are unlikely to arise by chance, even in isolation. A number of the classical sequence patterns for the superfamily in question should be present in the correct positions in the sequence in relation to other conserved patches and the conserved sequences should be consistent with a structural prediction for the relevant domain in cases where the domain structure is known. Thus there should be hydrophobic amino acids in the positions predicted to be in-pointing to stabilize the tertiary fold and the Cys residues should potentially be able to form disulphide bonds that are consistent with the fold. All these considerations were applied to the analysis of the IgSF relationship of the CD2 and CD4 antigens where there has been controversy concerning CD2 domain 1 and CD4 domain 2. In both cases it was shown by structure determination that these domains were in the IgSF and that correct predictions were made for the β strands in both domains ⁷⁻¹⁰.

A further analysis using the ALIGN program is illustrated in Table 3 for the cytokineR superfamily. This is one of the most diverse superfamilies in terms of sequence alignments as can be seen from Fig. 3. The ALIGN scores clearly support the superfamily relationship for the grouped sequences with the exception of the domains nominated for the IL4 and IL7 receptors. In the case of the IL4R domain only 2/10 scores are 3 SD or greater with another 4 scores being >2 SD. Thus the case for inclusion of the IL4R domain in the superfamily is weaker on the basis of the ALIGN scores. However, inspection of the conserved sequence patterns in Fig. 3 leaves little doubt that the IL4R domain is in the cytokineR superfamily. Cysteine residues can be confidently placed at all of the conserved positions and other conserved patterns are also present in the correct positions. For the IL7R domain the situation is much more ambiguous. The ALIGN scores are very weak and only a hint of the conserved sequence patterns is seen in the alignments. This domain would not be considered for inclusion in the cytokineR superfamily except that the domain is present in a cytokine receptor. The case for this will ultimately require validation by three-dimensional structure determination.

Domain sequence and structure: divergent and convergent evolution

In the above section criteria for defining a superfamily have been based on identifying a sequence pattern that is shared in a non-trivial way between sequences of different molecules. It is then argued that the presence of the sequence pattern indicates a relationship in evolution such that the domains that share the sequence pattern both derive from one original primordial domain. However, it could be argued that a certain structure dictates a sequence pattern and the sharing of the pattern is due to convergent evolution from different molecules rather than divergent evolution from a primordial domain. Consequently it may be found that sequences with no detectable common pattern form similar tertiary structures and thus that these are in the same superfamily even though there is no detectable sequence relationship.

It now seems very unlikely that a general structure will dictate a unique

sequence pattern. The rise to domains with no convincing fold. These are cytokineRSF proteins to a domain within the IgSF required to detect convergent evolution sequence patterns.

The converse of structure having no structure should be that sequences referred to seem useful. A number of small numerous occurrences of the Ig-fold may have been solution to them were not detected. There is no way to generate it seems best. This is sensitive data are multiple superfamilies. Superfamilies the sequence patterns in them.

Given that arises as to cell surface proteins the extracellular require diverse sequence patterns and usually out-pointing question arises structure of.

For cell surface the requirement. The small, they may have evolved coat protein the evolution cell differences enzymes and molecules

sequence pattern. This can be seen from a consideration of sequences that can give rise to domains with the Ig-fold. There are now five different sets of sequences with no convincing sequence similarity between them that can all give rise to the fold. These are the sequences of the Ig superfamily, the FN type III SF and cytokine RSF plus two sequences in the Pap-D bacterial protein that each give rise to a domain with an Ig-fold ¹¹⁻¹⁴. There is also enormous diversity of sequences within the IgSF that leads to the argument that there is no unique sequence required to determine any part of the IgSF-fold. Thus it seems rather unlikely that convergent evolution to yield the same structure would give rise to any common sequence pattern.

The converse argument is that all the sequences that give the same fold must have structure derived by divergent evolution and that all sequences with this structure should be included in the same superfamily. For example, the five sets of sequences referred to above might all be considered as IgSF sequences. It does not seem useful to take this point of view since there may be a relatively limited number of small stable protein folds that can occur and these may have evolved on numerous occasions in evolution. In this case each of the sets of sequences within the Ig-fold would have an independent primordial ancestor. Alternatively, there may have been one primordial structure which acquired mutations such that a solution to the structure was produced, ultimately giving rise to sequences that were not detectably similar to the ancestor family of sequences. At this stage there is no way to estimate the probability of the divergent versus the convergent route for generation of the same structure without recognizable sequence similarity and it seems best to stick to sequence patterns as the criteria for defining superfamilies. This is sensible from a practical as well as a theoretical standpoint since sequence data are much more readily obtained than tertiary structural data and the superfamilies defined on the basis of sequence would be grouped as subsets within superfamilies based on tertiary structure considerations. It seems better to retain the sequence criterion and to note that certain superfamilies have the same folding patterns in their domains.

Given that the same structure can arise from various sequences, the question arises as to why sequence patterns are conserved in evolution. Molecules on the cell surface present unique determinants for interaction with a soluble molecule, the extracellular matrix or with other cell surface receptors. Such interactions require diversity between molecules and not conservation of epitopes. The sequence patterns shared within a superfamily conserve the fold of the molecule and usually involve residues pointing inwards in the folded structure rather than out-pointing residues that are available for biological interactions. Thus the question arises as to what evolutionary force can operate to preserve the tertiary structure of the molecule?

For cell surface molecules it can be argued that the key evolutionary pressure is the requirement for molecular stability and, in particular, resistance to proteolysis. The small, tightly folded domains that make up most of the leucocyte molecules may have evolved as parts of stable coat proteins on single cell eukaryotes ^{6, 7}. These coat proteins then gave rise to the families of molecules that evolved along with the evolution of multicellular organisms, to mediate cell division and regulation of cell differentiation. Surface molecules are generally resistant to proteolysis by enzymes and this resistance is based on the folded structure, since denatured molecules are easily digested. One could argue that mutation to give new

recognition epitopes would be constrained by the necessity of preserving structure of the domain. In general this led to preservation of certain patterns that determine one particularly stable solution for the fold. alternative sequence patterns may exist that could also give a stable fold, but to reach these a number of simultaneous mutations may be required and a switch to a new pattern may be a rare event in evolution. If a new pattern this may become the founder of a new set of sequences in which the new pattern is retained, again because of the pressure of proteolysis. From this viewpoint it is likely that the Ig, FN type III and cytokineR superfamilies all arose from a common ancestor via sequence shifts as described above. This view might be favoured because domains of these superfamilies are found in molecules with similar functions and often a molecule may contain both Ig superfamily domains and domains of the FN type III and cytokineR superfamilies. In particular, Ig and FN type III domains are often found together in a single polypeptide.

Genomic structure and evolution of proteins with mixtures of domain types
The number of domains in a cell surface protein can vary greatly. In the case of the Thy-1 antigen there is a single IgSF domain making up the whole of the extracellular segment, whilst for the complement receptor 1 protein (CD35) the extracellular region consists of 30 CCPSF domains in a linear array. In these cases only one domain type is present but in other molecules there can be a mixture of domain types. For example the L-selectin (LECAM-1) antigen contains C-type, EGFSF and CCPSF domains.

The efficient build-up of proteins from individual domains during evolution appears to depend on two aspects of genomic structure. There should be an approximate concordance of the domain ends with intron/exon boundaries and the position of the intron with respect to the reading frame of a gene should be such that an open reading frame results from the recombination of an exon into the intron of an existing sequence¹⁵. Introns that are inserted after the first codon are called phase 1, those after the second base, phase 2, and those after the third base, phase 0¹⁶. Analysis of the intron/exon boundaries of domains present on leucocyte surface molecules shows that for most domain types each exon is of the same phase (usually 1) as illustrated in Table 4. Recombination of such exons will lead to the construction of new open reading frames. The domain does not need to be contained within a single exon to allow shuffling as long as the outermost intron boundaries are compatible. For instance, some IgSF domains are coded for by two exons² and the cytokineR domain in the IL2 receptor is coded for by three exons. In the latter case the internal splice sites are phase 2 and whilst the external ones are phase 1, thus it is not possible to get part of the domain integrated into a sequence containing phase 1 splice sites.

THE SUPERFAMILIES THAT ARE FOUND IN LEUCOCYTE CELL SURFACE MOLECULES

The superfamilies that are present in leucocyte surface molecules are discussed below together with alignments of some of the domains or repeat sequences. The alignments were made using a variety of computer programs (ALIGN¹⁷, PILEUP¹⁸) and then modified after visual examination. The ends of the

Table 4. Exon

Domain or re

Complement
CytokineR
EGFSF
Fibronectin
Fibronectin
IgSF
Lectin C-type
Lectin C-type
Lectin S-type
Leucine-rich
LinkSF
LDLRSF
Ly-6SF
MHC
Nerve growth
ScavengerR
Somatomedin

In both the Ig
two exons and
available on
numbers of
NK, not known

can be difficult
consideration
continues to
highly likely
domain 3 was
the domain
an asterisk
being a conserved
example, in
conserved C
system their
structural
the sequence

Table 4. Exon organization of domains and repeats of leucocyte surface molecules

Domain or repeat type	Do the domain boundaries coincide with introns with same splice sites?	Splice site	Usual number of exons per domain
Complement control protein (CCP)SF	Yes	type 1	1
CytokineRSF	Yes	type 1	2
EGFSF	Yes	type 1	1
Fibronectin type IISF	Yes	type 1	1
Fibronectin type IIISF	Yes	type 1	1
IgSF	Yes	type 1	1 or 2
Lectin C-typeSF (e.g. selectins)	Yes	type 1	1
Lectin C-typeSF (e.g. Kupffer cell receptor)	No	NA	3
Lectin S-typeSF	NK	NK	NK
Leucine-rich glycoprotein repeat	No	NA	NA
LinkSF	Yes	type 1	1
LDLRSF	Yes	type 1	1
Ly-6SF	No	NA	NA
MHC	Yes	type 1	1
Nerve growth factor receptor (NGFR)SF	No	NA	NA
ScavengerRSF	NK	NK	NK
Somatomedin BSF	Yes	type 1	1

In both the IgSF and the CCPSF domains there are examples where the domain is encoded by two exons and also where two domains are encoded by one exon. Only limited data are available on some of the domains and it is possible that other examples with different numbers of exons per domain or motif may be found.
NK, not known, NA, not applicable.

can be difficult to define from the sequence and this problem is illustrated by consideration of the structure for CD4. In CD4 the last β strand of domain 1 continues directly into domain 2 and between CD4 domains 3 and 4 it seems highly likely that the overlap will be even greater and that the last β strand of domain 3 will also be the first β strand of domain 4. Thus in the alignments shown, the domains are defined with respect to key internal residues that are marked with an asterisk, and the beginnings and ends can be taken for statistical comparison as being a constant number of residues before and after the conserved position. For example, in the case of the IgSF this is taken as 20 residues before and after the conserved Cys positions. If the goal was to express a single domain in an expression system then sequence alignments and structure should be taken into account and a structural prediction would be attempted on the basis of all the data to decide on the sequence that should be expressed.

The complement control protein (CCP) superfamily (Figs 1 and 2)

This domain is named CCP because it is commonly found in proteins that control the complement cascade ¹⁹. For instance, factor H consists solely of 20 CCPSF domains whilst other complement components contain CCPSF domains mixed with other domains, e.g. factors B and C2 each contain three CCPSF domains together with a serine protease domain. The CCP domain is also commonly called the short consensus repeat or SCR ¹⁹. It is present in widely different numbers in cell surface molecules ranging from 30 domains in complement receptor 1 (CD35) to a single domain in L-selectin. These domains are clearly involved in protein binding and the CR1 (CD35) and CR2 (CD21) complement binding regions have been mapped to the first four CCP domains of each of the first three groups of seven domains in CD35 and to the first two domains of CD21.

The structure of one CCPSF domain from complement control protein factor H has recently been solved using NMR and consists of two segments of antiparallel β sheet and a short triple-stranded β sheet with no α -helical structure ²⁰. The folding pattern for this domain is shown in Fig. 2 and the β strand positions are marked above the sequence alignments shown in Fig. 1.

Cytokine receptor (cytokineR) superfamily (Figs 3 and 4)

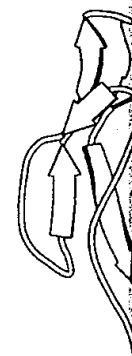
Three domain types are found in cytokine receptors including those of the Ig, FN type III and cytokineR superfamilies. A common arrangement is to have a single NH₂-terminal cytokineRSF domain followed by an FN type IIISF domain, but there are variations on this theme. Initially these two domain types were not distinguished ⁵ and the term haematopoietin receptor superfamily was widely used for molecules containing this pair of domain types ^{21,22}. We use the term cytokine receptor superfamily for the domain of about 100 amino acids usually found NH₂-terminal to the FN type IIISF domain and alignments of domains from this superfamily are shown in Fig. 3. Analysis with the ALIGN program (Table 3) gives good evidence for the presence of the cytokineRSF domain in the receptors for IL2 (β chain), IL3, IL6, growth hormone, granulocyte-macrophage colony stimulating factor, erythropoietin and in the GMP130 protein ²². The presence of a cytokineRSF domain in IL4R is less strongly supported by ALIGN analysis but as discussed above the case for inclusion of this domain is convincing if all the data are considered. The IL7 receptor contains a clear FN type IIISF domain but the sequence at the NH₂-terminal region shows only a distant relationship to the cytokineRSF domains [see p.338]. The possible cytokineR domain in the IL7 receptor ^{22,23} is shown below the other sequences in Fig. 3 but the correctness or otherwise of this assignment will require validation by tertiary structure determination.

The structure of the growth hormone receptor has recently been solved by X-ray crystallography ¹³ and this has revealed the fold for the cytokineRSF and the FN type IIISF domains that constitute the extracellular domain of this receptor (an FN type IIISF domain has also been solved by NMR - see below). These domains have similar folds that are also similar to the folds of IgSF C2 set domains ^{8,9} and the PapD chaperone protein domains ¹¹. Bazan ²⁴ had previously argued that there may be structural similarities between cytokineRSF domains, FN type IIISF domains and IgSF domains on the basis of predicting patterns of β strands in the sequences. Despite the success of these predictions the degree of sequence similarity between these domain types is low. The cytokineRSF domains have a characteristic Cys-X-Tip sequence together with three other conserved Cys residues, whilst the FN type

Factor H	L P C K S - P P E
CD35d12	R V C Q P - P P D
Factor B	G S C S L - - - E
L-Selectin	I Q C E - P L E
C4BPA	N S C I N - L P D
IL2R1	E L C D D D P P E
FXIII	E P C T V - N V D

Factor H	F G I D G P A -
CD35d12	Y D L R G A A -
Factor B	F Y P Y P V Q -
L-Selectin	T N L T G I E -
C4BPA	Y K P T T D E P
IL2R1	F R R I K S G S I
FXIII	Y D L S P L T P I

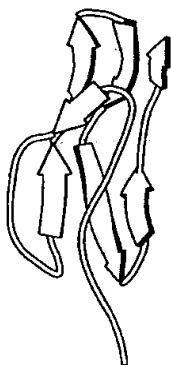
Figure 1. CCP superfamily are boxed. The lines above strands determined from the (see Fig. 2) ²⁰. The asterisks shown on the figures to identify the entries in Section II. The Swissprot database unless a number and residue number factor H precursor domain 1 precursor domain 12 (P17 factor B (P00751, 10-74); L-C4BPA, complement C4-binding receptor α chain precursor (P05160, 452-516)



Factor H CCPSF

Factor H	L P C K S - P P E I S H G V V A H M - - - - S D S Y Q Y G E E V T Y I	F E G
CD35d12	R V C Q P - P P D V L H A E R T Q R D K - - - - D N F S P G Q E V F Y I	E P G
Factor B	G S C S L - - - E G V E I K G G S F R L - - - - L Q E G Q A L E Y I	P S G
L-Selectin	I Q C E - - P L E A P E L G T M D C T H P - - L G N F N F N S Q C A F I	S E G
C4BPA	N S C I N - L P D I P H A S W E T Y P R P T K E D V Y V V G T V L R Y I	H P G
IL2R1	E L C D D D P P E I P H A T F K A M - - - - A Y K E G T M L N C I	K R G
FXIII	E P C T V - N V D Y M N R R N I E M K W - K Y E G K V L H G D L I D F I	K Q G
Factor H	F G I D G P A - - - - I A K C L G - E K W S H P - - - - -	C I
CD35d12	Y D L R G A A - - - - S M R C T P Q G D W S P A A - - - - -	C E
Factor B	F Y P Y P V Q - - - - T R T C R S T G S W S T L K T Q D Q K T V R K I	C R
L-Selectin	T N L T G I E - - - - E T T C E P F G N W S S P E - - - - -	C Q
C4BPA	Y K P T T D E P T - - - - T V I C Q K N L R W T P Y Q - - - - -	C E
IL2R1	F R R I K S G S L - - - - Y M L C T G N S S H S S W D N Q C - - - - -	C T
FXIII	Y D L S P L T P L S E L S V Q C N R - G E V K Y - - - - -	C T

Figure 1. CCP superfamily domains. Residues identical in four or more sequences are boxed. The lines above the sequences correspond to the positions of the β strands determined from the structure of factor H domain 16, residues 92-985 (see Fig. 2) ²⁰. The asterisks mark the positions of the conserved residues shown on the Figures to identify domains in each entry for a molecule as the entries in Section II. The sequences of the following proteins are from the Swissprot database unless otherwise indicated and the database accession number and residue numbers are given in brackets. Factor H, human complement factor H precursor domain 16 (P08603, 929-985); CD35d12, complement factor H precursor domain 12 (P17927, 745-799); Factor B, HR16 human complement factor B (P00751, 10-74); L-selectin, L-selectin precursor (P14151, 195-252); C4BPA, complement C4-binding protein (P04003, 249-313); IL2R1, interleukin-2 receptor α chain precursor (P01589, 22-83); FXIII, coagulation factor XIII chain precursor (P05160, 452-516).



Factor H CCPSF domain

Figure 2. The folding pattern of the CCPSF domain. Ribbon diagram showing the folding pattern of a CCPSF domain from factor H determined by NMR ²⁰. The β strands are shown as broad arrows pointing from the amino to carboxy direction and the connecting loops as thinner lines.

Figure 3. Cytokine receptor superfamily

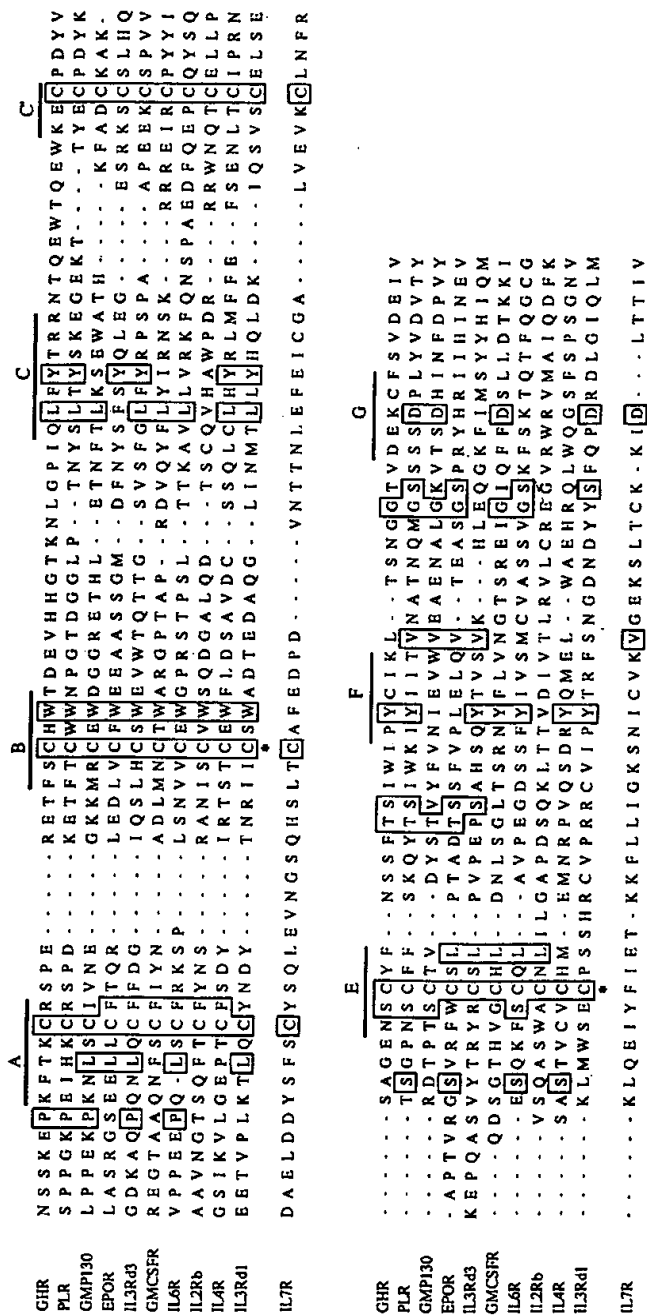
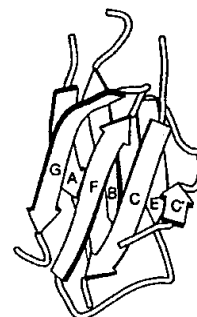
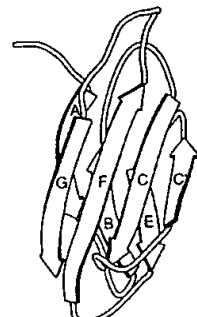


Figure 3. (opposite) Cytokine receptor superfamily. More sequences are boxed. Residues that are shown on a molecule as shown in Section 3.1.1. the Swissprot database unless otherwise indicated. IL3 receptor precursor (P10912, 21-116); GMPI30, human IL3 receptor precursor (P10912, 124-218); EPOR, mouse erythropoietin receptor precursor (P10912, 243-347); GMCSFR, human granulocyte-macrophage colony-stimulating factor receptor precursor (P14784, 26-125); IL7R, human IL7 receptor precursor (P14784, 26-125). The start corresponds to residue 20 amino acids. Not more difficult to define due to the close proximity of the predicted boundaries of the type IIIISF domains in GHR.



Human GHR domain 1



Fibronectin domain 21

Figure 3. [opposite] CytokineR superfamily domains. Residues identical in four or more sequences are boxed. The asterisks mark the positions of the conserved residues that are shown on the figures to identify domains in each entry for a molecule as shown in Section II. The sequences of the following proteins are from the Swissprot database unless otherwise indicated and the database accession number and residue numbers are given in brackets. GHR, human growth hormone receptor precursor (P10912, 46–147); PLR, rat prolactin receptor precursor (P05712, 21–116); GMP130, human membrane glycoprotein gp130 precursor (PIR: A36337, 124–218); EPOR, mouse erythropoietin receptor precursor (P14753, 42–140); IL31, mouse IL3 receptor precursor domains 1 and 3 (PIR: A35782, d1, 29–127; d3, 243–347); GMCSFR, human GM-CSFR precursor (P15509, 116–214); IL6R, human IL6 receptor precursor (P08887, 112–214); IL2R β , human IL2 receptor β chain precursor (P14784, 26–125); IL4R, mouse IL4 receptor precursor (P16382, 24–122); IL7R, human IL7 receptor precursor (P16871, 32–127). The sequence alignments are from 20 amino acids NH₂-terminal from the conserved CXW. The sequence start corresponds to residue 2 in the prolactin receptor. The COOH-terminus is more difficult to define due to the lack of conserved residues and that shown is close to the predicted boundary between the cytokineRSF domains and the FN type IIISF domains in GHR, PLR, IL6R.

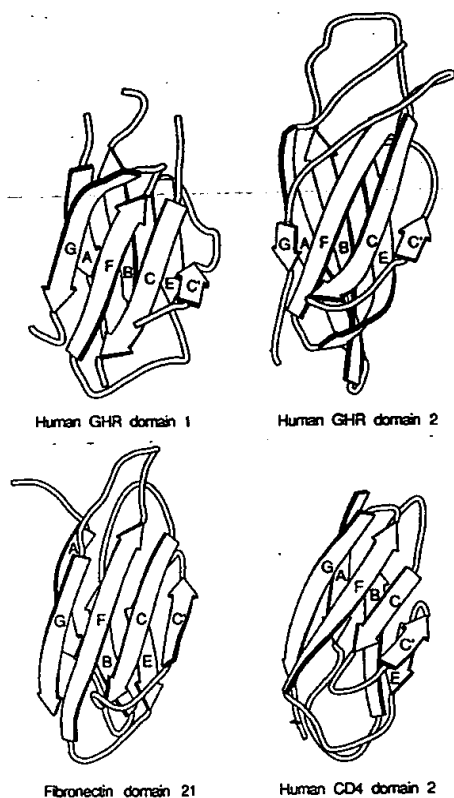


Figure 4. The folding pattern of the cytokineR superfamily and fibronectin type IIISF domains. Ribbon diagrams of the cytokineR superfamily and FN type IIISF domains from human growth hormone receptor¹³, and FN type IIISF domain 21 from human fibronectin¹⁴. The IgSF C2-s domain 2 from human CD4 is included for comparison^{8,9}. The β strands are shown as broad arrows pointing from the amino to carboxy direction and the connecting loops as thinner lines. Some gaps are present in the loops of the growth hormone receptor where the structure has not been fully resolved¹³. Each β strand is labelled using the same nomenclature as in the IgSF. This lettering corresponds to that in the sequence alignments (Figs 3, 8, 12).

IIISF domains lack a conserved pattern of Cys residues. Using the ALIGN program 100 comparisons between cytokineRSF and FN type IIISF domains were made and of these 81 gave a score less than 2 SD and only eight comparisons gave scores greater than 3.0 SD (5.9, 5.6, 4.2, 3.7, 3.5, 3.4, 3.3, 3.0 SD; unpublished observations). Although there were two good scores out of 100 and some moderate ones, the case for a relationship in evolution based on ALIGN analysis is much weaker than that for members of the individual superfamilies. The possibility of an origin by divergent evolution for the cytokineR, FN type IIISF and IgSF domains is discussed above.

Several members of the cytokineR superfamily show small patches of sequence similarity in their cytoplasmic domains. These are reviewed in ref. 25.

Epidermal growth factor (EGF) superfamily (Figs 5 and 6)

EGFSF domains are found in EGF itself and in transforming growth factor (TGF) α . This domain is also found in a variety of secreted proteins such as blood coagulation factor IX and cell surface molecules such as in the selectins L-selectin, E-selectin and P-selectin (CD62). The structures of EGF, TGF α and the factor IX EGFSF domain have recently been determined and show similarity in folding pattern²⁶⁻²⁸. The structures of EGF and factor IX EGFSF domains are shown in Fig. 6. The latter is slightly smaller than EGF itself but is probably representative of the repeating EGFSF domains found in many proteins (see Fig. 5). The single EGFSF domain from factor IX has functional activities distinct from the EGF itself, for example it has Ca^{2+} -binding activity²⁸. It is likely that EGFSF domains are

FA9-1	V	D	G	D	Q	C	-	-	E	S	N	P	C	L	N	G	S	C	K	D	-	D	I	N	S	Y	E	C	W	C	P	F	G	F	E	G	K	-	-	N	C	E	L			
FA9-2	-	-	D	V	T	C	N	I	-	K	N	G	R	C	-	E	Q	F	C	K	N	S	-	A	D	N	K	V	V	C	S	C	T	E	G	Y	R	L	A	E	N	K	S	C	E	P
EGF	N	S	D	S	E	C	P	L	S	H	D	Q	Y	C	L	H	D	G	V	C	M	Y	I	E	A	L	D	K	Y	A	C	N	C	V	V	G	Y	I	G	E	-	-	R	C	Q	Y
L-Sel	-	-	T	A	S	C	-	-	Q	P	W	S	C	S	O	H	G	E	C	V	E	-	T	I	N	N	Y	T	C	N	C	D	V	G	Y	Y	G	P	-	-	Q	C	Q	F		
CD62	-	-	T	A	S	C	-	-	Q	D	M	S	C	S	K	Q	G	E	C	L	E	-	T	I	G	N	Y	T	C	S	C	Y	P	G	F	Y	G	P	-	-	E	C	E	Y		
E-Sel	-	-	T	A	A	C	-	-	T	N	T	S	C	S	G	H	G	E	C	V	E	-	T	I	N	N	Y	T	C	K	C	D	P	G	F	S	Q	L	-	-	K	C	E	Q		
PRTC	P	L	E	H	P	C	-	-	A	S	L	C	C	G	H	G	T	C	I	D	-	G	I	G	S	F	S	C	D	C	R	S	G	W	E	G	R	-	-	F	C	Q	R			
114/A10	G	P	S	D	L	C	-	-	N	P	N	P	C	X	G	T	A	S	C	V	K	-	L	H	S	K	H	F	C	L	C	L	E	G	Y	Y	N	S	S	L	S	S	C	V	K	
NOTCH	T	N	D	E	D	C	-	-	T	E	S	S	C	L	N	G	S	C	I	D	-	G	I	N	G	Y	N	C	S	C	L	A	Q	Y	S	G	A	-	-	N	C	Q	Y			

Figure 5. EGF superfamily domains. Residues identical in five or more sequences are boxed. The asterisks mark the positions of the conserved residues that are marked on the domain organization figures in the entries in Section II. The sequences of the following proteins are from the Swissprot database and the database accession number and residue numbers are given in brackets. FA9-1 and FA9-2, human coagulation factor IX precursor (P00740, 92-130 and 131-172); EGF, human epidermal growth factor precursor (P01133, 971-1014); L-Sel, human L-selectin precursor (P14151, 157-193); CD62, human CD62 or P-selectin precursor (P16109, 160-196); E-Sel, E-selectin precursor (P16581, 138-176); PRTC, human protein C precursor (P04070, 96-133); 114/A10, mouse haematopoietic cell surface protein 114/A10 precursor (P19467, 232-274); NOTCH, Drosophila notch protein (P07207, 1021-1059). The ends of the alignment correspond to those of the coagulation factor IX EGFSF domain whose structure has been determined²⁸. The structure of EGF itself has been determined for a sequence that extends a further four residues beyond that shown (see Fig. 6)²⁷. The bars above the sequence indicate the positions of the β -strands.



EGF

Fibr d7	T	A	V	T
Fibr d8	T	V	L	V
Mannose R	Y	E	A	M
Factor XII	K	A	E	E
Collag	R	A	D	S

C	T	I	E	G	R	Q	D	G
C	T	S	E	G	R	R	D	N
C	T	S	A	G	R	S	D	G
C	T	H	K	G	R	P	G	P
C	T	S	E	G	R	G	D	G

Figure 7. Fibronectin type II sequences are boxed. The sequences of the following proteins are from the database accession number and residue numbers are given in brackets. human fibronectin precursor (P02751, 1-246); human mannose receptor (P15111, 1-100); human coagulation factor XII precursor (P00740, 131-172); collagenase precursor (P00740, 131-172). The sequences are based on the exon b

recognition structures of the 36 EGFSF domain necessary for the interaction

Fibronectin (FN) type II
The FN type II domain sequence patterns with

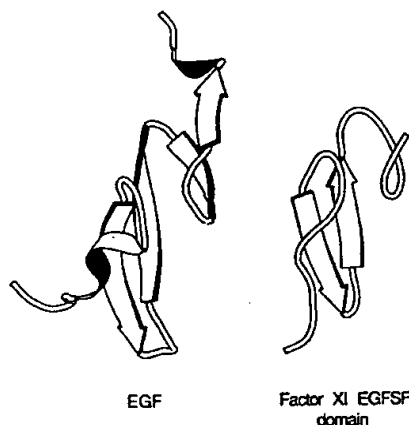


Figure 6. The folding pattern of EGFSF domains. Ribbon diagrams of EGF²⁷ and a coagulation factor IX EGFSF domain²⁸. The β strands are shown as broad arrows pointing from the amino to carboxy direction and the connecting loops as thinner lines. The NH_2 -terminal core of the structure is similar in both domains but the EGF structure extends further with two more short β strands.

Fibr d7	T	A	V	T	Q	T	Y	G	G	N	S	N	G	E	P	C	V	L	P	P	T	Y	N	G	R	T	F	Y	S
Fibr d8	T	V	L	V	Q	T	Q	G	G	N	S	N	G	A	L	C	H	F	P	F	L	Y	N	N	H	N	Y	T	D
Mannose R	Y	E	A	M	Y	T	L	L	G	N	A	N	G	A	T	C	A	F	P	F	K	F	E	N	K	W	Y	A	D
Factor XII	K	A	E	E	H	T	V	V	L	T	V	T	G	E	P	C	H	F	P	F	Q	Y	H	R	Q	L	Y	H	K
Collag	R	A	D	S	T	V	M	G	G	N	S	A	G	E	L	C	V	F	P	F	T	F	L	G	K	E	Y	S	T

C	T	T	E	G	R	Q	D	G	H	L	W	C	S	T	T	S	N	Y	E	Q	D	Q	K	Y	S	F	C	T	D	H	T
C	T	S	E	G	R	R	D	N	M	K	W	C	G	T	T	Q	N	Y	D	A	D	Q	K	F	G	F	C	P	M	A	A
C	T	S	A	G	R	S	D	G	W	L	W	C	G	T	T	T	D	Y	D	T	D	K	L	F	G	Y	C	P	L	K	F
C	T	H	K	G	R	P	G	P	Q	F	W	C	A	T	T	P	N	F	D	Q	D	Q	R	W	G	Y	C	L	E	P	K
C	T	S	E	G	R	G	D	G	R	L	W	C	A	T	T	S	N	F	D	S	D	K	K	W	G	F	C	P	D	Q	G

Figure 7. Fibronectin type IISF domains. Residues identical in three or more sequences are boxed. The asterisks mark the positions of the conserved residues that are marked on the domain organization figures in the entries in Section II. The sequences of the following proteins are from the Swissprot database and the database accession number and residue numbers are given in brackets. Fibr, human fibronectin precursor (P02751, d7, 314–373; d8, 374–434); Mannose R, human mannose receptor precursor (P22897, 153–212); Factor XII, human coagulation factor XII precursor (P00748, 32–91); Collag, human type V collagenase precursor (EC 3.4.24.7) (P14780, 332–391). The ends of the alignment are based on the exon boundaries of the fibronectin domains.

recognition structures that can be involved in various functions. For example, two of the 36 EGFSF domains of the *Drosophila* protein *Notch* have been shown to be necessary for the interaction of *Notch* with the *Delta* and *Serrate* proteins²⁹.

Fibronectin (FN) type II superfamily (Fig. 7)

The FN type II domains were first identified as one of three different repeating sequence patterns within the fibronectin molecule. The FN type IISF domain has

been found in few other proteins and the only leucocyte molecule with this domain is the mannose receptor which contains one FN type IIISF domain. The structure of a sequence from bovine seminal fluid protein PDC-109 that shows sequence similarity over part of the FN type II domain alignment shown in Fig. 7, has been determined by NMR ³⁰.

Fibronectin (FN) type III superfamily (Figs 4 and 8)

The FN type IIIISF domain was identified in an extracellular matrix protein but is also common in membrane molecules and particularly these found in the nervous system which often have IgSF domains ². It has also been found in large numbers in the group of muscle proteins that bind myosin such as twitchin in *Caenorhabditis elegans* ³¹ and also in titin in mammals ³². This is the only group of IgSF molecules found so far in the cytosol. Another example of cytosolic localization of FN type IIIISF domains is in the cytoplasmic segment of the integrin β_4 chain ³³ (note the external regions of integrins do not contain any FN type IIIISF domains). This is currently the only example of a domain found at the surface of leucocytes which is also present on the cytoplasmic side of a transmembrane protein.

Structures for FN type IIIISF domains have recently been solved by NMR ¹⁴ and X-ray crystallography ¹³. This domain consists of two β sheets with a similar folding pattern to the IgSF fold, the CytokineR domain and the domains of the PapD chaperone protein ¹¹. However, there is no significant sequence similarity amongst these proteins as analysed by the methods discussed above.

Immunoglobulin (Ig) superfamily (Figs 9–12)

The immunoglobulin superfamily (IgSF) is the largest superfamily of cell surface proteins in general and for leucocyte antigens in particular, as is evident from the collated data in Table 1 in Chapter 1 which shows that approximately 40% of leucocyte membrane polypeptides contain IgSF domains. The structures of several IgSF domains have been determined by X-ray crystallography including Ig V- and

Figure 8. (opposite) *Fibronectin type IIIISF domains. Residues identical in five or more sequences are boxed. The positions of the β strands determined for domain 21 of human fibronectin are indicated above the sequences ¹⁴. See Fig. 4 for folding patterns of FN type IIIISF domains from fibronectin and growth hormone receptor. The asterisks mark the positions of the conserved residues that are marked on the domain organization figures in the entries in Section II. The sequences of the following proteins are from the Swissprot database unless otherwise indicated and the database accession number and residue numbers are given in brackets. GHR, human growth hormone receptor precursor (P10912, 148–251); FIBR, human fibronectin precursor (P02751:d12 605–700, d13 719–800; d16 996–1085, d21 1447–1541); LAR, human LAR precursor (P10586, 596–694); TWIT, twitchin cytoplasmic protein from *Caenorhabditis elegans* (PIR:S07571 1761–1854); CAML1, mouse neural adhesion molecule L1 precursor (P11627, 916–1012); IL7R, human interleukin 7 receptor precursor (P16871, 128–231); GMP130, human membrane glycoprotein gp130 precursor (PIR: A36337, 221–32); PLR, rat prolactin receptor precursor (P05710, 121–224); IL3LR, mouse IL3-receptor-like protein precursor (AIC2B) (PIR: A35782, d2; 135–243, d4; 342–441)*

Figure 8. Fibronectin type IIIISF domains

	A	B	C
FIBR-21	VSDVPRDLEVAATPT	SLLSWDAPAVT	VRYRIYGETG
GHR	QPDPIALNWTLLNVSLTGIHADIQYRWAEAPRNADIQK	GWMLVLEYELQYKEVN	
FIBR-12	YPSSSGPVEVETETPS	QPNSHPIQWNAIPQPSH	ISKYLIRWRPKN
FIBR-13	SPLVATSESVTEITA	SSSFVSWVSAQSDT	VSGFRVEYELS
FIBR-16	KLDAPTNQEFVNET	DSTVLVRWTPPRAQ	ITGYRLTVGLTR

Figure 8. Fibronection type IIISF domains

	A	B	C		C'	E	F	G
FIBR-21	VSDVPRDLT <u>EV</u> AAATPT	<u>S</u> LLISWDA <u>PA</u> VT	...VRYRIT <u>Y</u> GETG		GNSPVQEFVTVPGS	...KSTATISGLK <u>PGVD</u> YTI <u>TV</u> AVTGR <u>G</u> DS	...ASKPI SINRTE	
GHR	QDPPIALNWTLLNVSLTG	AD <u>Q</u> VRWEAPRNADIQK	...GWMVLE <u>Y</u> ELQYKEVN		<u>E</u> TKWKMDPIIL	...TT <u>S</u> VPVYSKVDKEYEVRVRSKQRNSG	...NYGEFSEVLYVT	
FIBR-12	YPSSSGPVEVFIETPS	QPN <u>S</u> HPIQWNAPOPSH	... <u>L</u> SKYLILRWPKN		SVGRWKEATIPG	...HLNSYTIKGLKPGVYEGQLISIQYCH	...QEVTRPFDFTTT	
FIBR-13	SPLVATSEVTEITA	<u>S</u> SFVWSVSAASDT	...VSGFRVE <u>Y</u> ELS		<u>E</u> GDEPQYLDLPS	...TATSVNIPDLPPGRKYIVNVYQISEDGE	...QSLILSTSQTT	
FIBR-16	KLDAPTNLQFVNET	DSTLVLRWTPPRAQ	...ITGRLTVGLTR		RGQPRQYNVGPS	...VSKYPLRNLQPASEYTVSLVAIKGNOE	...SPKATGVFTTL	
LAR	PSAPPQKVMCVSMGS	TTVRVSWVPPADSRNG	...VITQYSAHEAVD		GDRGRHVVDGIS	REHSWDLVGLKXWTEYRVRVAHTDVCPG	...PESSPVLVRID	
TWIT	RPDRPPGRPEPTDWS	DHVDLKWDPPLSDGGA	...PIEEYQIEKRTKY		GRWEPI TVPG	GQTTATVPDLTPNEEYEFRVAVNKGGP	...SDPSDASKAVI	
CAML1	VPGHPEALHLECQSD	<u>T</u> SL <u>L</u> HMQPPLSHNG	...VLTCYLLSYHPVE		GESKEQLFFNLSDP	ELRTHNLTNLNPDLYRQFQLQATTQGG	...PGIAIVREGGIM	
IL7R	KPEAPFDLSIYREGA	NDFVVTENTSHLQKKYV	...KVLMDVAAYRQEK		DEKWTHTVNLSST	KLTLQRLQPAAMEYIKVRSIPDHYF	...KGFWSWSPSYFFRTPEI	
GMP130	KPNPPHNLSSVINSEELS	SILKLTTNP'SIKSV	...ILKYNIQYRTKD		ASTWSQIPPEDTASTRS	SFTVQDLKPFTEYVFRIRCMKEDG	...KGYMSDWSEASGITYED	
PLR	IVEPPRNLTLEVQQLKD	KKTYLWVKWSPPTITDVKT	...GWFTMEYEIRLKPBE		AEEWEIHFTG	HQTFKVFVDLYPGQRYLYQTRCKPDHG	...YMSRWSSQESSVEMPND	
IL3LRd2	QPLPKNYSISSSED	RFLLEWVSLSGDAQVSWLSSKDI	EEFEVAYKRLQ		DSMEDAYS LHTSKFQVNF	PEPKLFLPNSIYAPRVTRLYPGLSS	SSGRPSRWSPSAHWDSQPQ	
IL3LRd4	IQMEPPTLNLTKNRD	<u>S</u> YSLHMETQKMAV	...SFIEHTFQVYKKKS		DSWEDSKTENL	DRAH <u>S</u> MDLSQEPDTSYCARVRVKPISNY	...DGIWSEYTWKTDWY	
		</						

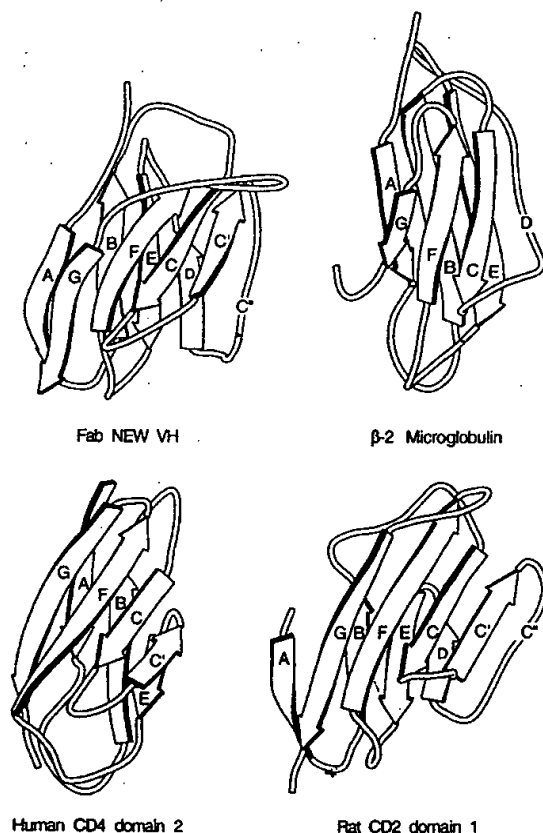


Figure 9. The folding pattern of IgSF domains. Ribbon diagrams for IgSF domains. Ig V set (V_H of human NEW Fab); Ig C1 set (β 2-microglobulin); Ig C2 set (CD4 domain 2); and Ig V set lacking the normally conserved disulphide between β strands B and F (rat CD2 domain 1). The β strands are shown as broad arrows pointing from the amino to carboxy direction and the connecting loops as thinner lines. These are labelled with the corresponding strand letters used in the alignments for the Ig V-set, C1-set and C2-set sequences (Figs 10–12) and in the FN type III ISF domains (Figs 4 and 8). The data are from the Brookhaven Protein Structure Database apart from CD2 ¹⁰.

Figure 10. (opposite) Immunoglobulin V-set domains. Residues identical in five or more sequences are boxed. The positions of the β strands are indicated above the sequences. The asterisks mark the positions of the conserved residues that are marked on the domain organization figures in the entries in Section II. The sequences of the following proteins are from the Swissprot database and the database accession number and residue numbers are given in brackets. Ig lambda, mouse Ig λ chain precursor (MOPC 104E) (P01724, 21–129); Ig kappa, human Ig κ chain Roy (P01608, 3–107); Ig heavy, human Ig heavy chain NEWM (P01825, 3–116); TcR beta, human TcR β chain precursor (P01733, 22–135); TcR alpha, mouse TcR α chain precursor (P01739, 23–132); CD8 beta, rat CD8 β chain precursor (P05541, 21–134); CD8 alpha, rat CD8 α chain precursor (P07725, 27–138); CD4d1, human CD4 precursor domain 1 (P01730, 21–123); Thy-1, rat Thy-1 precursor (P01830, 18–128); CD2 d1, rat CD2 precursor domain 1 (P08921, 20–120).

Figure 10. Immunoglobulin V-set domains

Protein

	A	B	C	C'
Ig lambda	AVVTQESALTTSPGELVLTCTRSSTGAV	TTSNYANVQKPP--DHLFTGLIGG--	TNNRPGV...	
Ig kappa	QMTQSPLSLSASVGRVITITCQASQD...	ISIFLNWYQKPG-KAPKLLIYDA...	SKLEAGV...	

Figure 10. Immunoglobulin V-set domains

	A	B	C	C'	C''
Ig lambda	AVVTQESALTTSPGGETVLTCTCRSSITGAVTTSNAYANWVQQKP--DHLFTGLIGG--TNNRAPGV---				
Ig kappa	QMTQSPSSSLASASVGDRTVTTCQASQD---ISIFLNWYQQKPG-KAPKLLIYDA---SKLEAGV---				
IgG heavy	GLEQSGPQLVVRPSQT-LSLTCTVSGS--TFSNDYYTAVRQPPG-RGLEWIGYVFYH--GTSDDTTPL				
TcR beta	GVIOQSPRHEVTEMGQEVTLRCKPISGH---NSLFWYRQTM-RGLELLIYFN--NNVPIDDSGMP				
TcR alpha	NVQQSPESLIVPEGARTSLNCTFSDS---ASQYFWYRQHSG-KAPKALMSIFS--NOEKE---				
CD8 beta	ALLQTPSSSLLVQTNQTAKMSCCAKTFPK---GTTIYWLRELQDSNKNKHFELASRTSTKGIKY---				
CD8 alpha	QLQLSPKKVDAEIQGEVKTCEVLRDTS---QGCWLFERNSSSELLQPTFIIVYSSRSKLNLDLD				
CD4 dI	AATQGGKKVVLGKKGGDTVELTCTASQKK---SIQPHMKNSN---QIKILGNQG---SFLTCKGPSKL				
Thy-1	RGORVISLTACLVNQNLRLDCRHEHNTNLPQHESLTRE---KKKHVLSGTL---GVPEHTY---				
CD2 dI	ADCRDSGTVWGALCHGININIPNFQMTD--DIDEVRWERGS---TLVAEFKRKMKPFLKSG----				
	D	E	F	G	
Ig lambda	PARFSGSLI---GNKAAITTTGAQTETDEAIFYFCALWYSN---HWVFGGGTKLTVL				
Ig kappa	PSRFSGTGS---GTDFTFTLSSLPEDIAITYYCQQFDNL---PLTEGGGKVDFK				
IgG heavy	RSRVTMLVDTSS--KNQFSRLSSVTAAADTAVYCCARNLIAG---C1DVWGGQSLVTVS				
TcR beta	EDRFSAKMPNA---SFSTLKTQPSEPRDSAVYFCASSFSTCSANYGYTFGSGTRITVV				
TcR alpha	EGREFIHLNKA--SLHFSLHIRDSQPSDSALYLCAVTLYOG-SGNKLIIFGTGTLTLLSVK				
CD8 beta	GERVKKNMTLSFNSITLPLKLMVDVKPEDSGYFCAMVGSF---MVVFGTGTKLTVV				
CD8 alpha	PNLFSAKKE---NNKYILTSLKFSSTKNQGYFCSITSNS---VMYFSPLPVPYFQK				
CD4 dI	NDRADSRRLWD-QGNFPLIKNLKIEDSDTYIICEVE---DQKEEVQLLVF				
Thy-1	RSRVNLFSDR---FIKVLTLANFTTKDEGDYMCDELRVSGQNPTSSNKTINIRDKLV				
CD2 dI	--AFEILA-----NGDLKIKNLTRDSGTINVTYSTNG---TRILDKALDRL				

Figure 11. Immunoglobulin C1-set domains

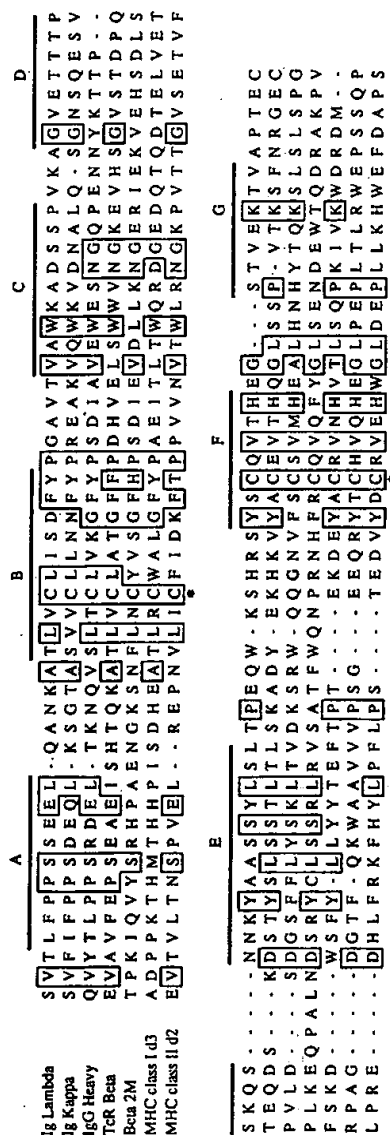


Figure 12. Immunoglobulin C2-set domains

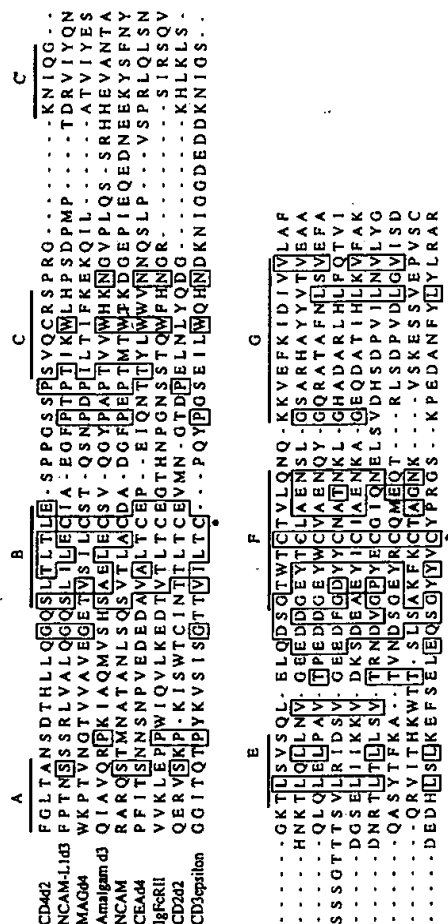


Figure 11. (opposite) Immunoglobulin C1-set domains. The sequences are boxed. The asterisks mark the sequences. The asterisks marked on the domain organization of the following proteins are otherwise indicated and the domain organization is given in brackets. Ig Lambda, Kappa, human Ig κ chain C region (P01857, 230-329); TcR β chain (P01857, 230-329); 2M, human β 2-microglobulin (P01857, 230-329); MHC Class I HLA α chain protein (P01857, 230-329); MHC Class II D (P01857, 230-329).

Figure 12. (opposite) Immunoglobulin C2-set domains. The sequences are boxed. The asterisks mark the sequences. The asterisks marked on the domain organization of the following proteins are otherwise indicated and the domain organization is given in brackets. CD4d2, human CD4d2 (P11627, 243-331); MAGd4, mouse neural cell adhesion molecule (P07722, 327-412); Amalgam d3 (P15364, 231-327); NCAM, neural cell adhesion molecule (P13590, 203-295); CEAd4, human CEAd4 (P06731, 321-410); IgFcRII, human IgFcRII (P06731, 321-410); CD2d2, human CD2d2 (P06731, 321-410); human CD3e precursor (P07722, 327-412).

domains, β 2-microglobulin 12 and 2, 8, 9 and CD8 α 35. The determined by NMR 10. These by sequence similarities over with distinct folding patterns fold consists of a sandwich of 5-10 amino acids with a but not all domains. The sequence in-pointing residues in the β connect the strands and the core of the fold is made up positioning of these is shown vary considerably in length being the archetype for the 1 domains forms an additional connection between these in antibody and TcR V-domain.

Figure 11. (opposite) Immunoglobulin C1-set domains. Residues identical in four or more sequences are boxed. The positions of the β strands are indicated above the sequences. The asterisks mark the positions of the conserved residues that are marked on the domain organization figures in the entries in Section II. The sequences of the following proteins are from the Swissprot database unless otherwise indicated and the database accession number and residue numbers are given in brackets. Ig Lambda, human Ig λ chain C region (P01842, 7-104); Ig Kappa, human Ig κ chain C region (P01834, 6-106); IgG Heavy, human Ig γ -1 C region (P01857, 230-329); TcR Beta, human TcR β chain (P01850, 10-117); Beta 2M, human β 2-microglobulin precursor (P01884, 24-119); MHC class I d3, human MHC Class I HLA α chain precursor domain 3 (PIR:A02189, 203-301); MHC class II d2, human MHC Class II DR α chain precursor domain 2 (PIR:A02206, 113-209).

Figure 12. (opposite) Immunoglobulin C2-set domains. Residues identical in four or more sequences are boxed. The positions of the β strands are indicated above the sequences. The asterisks mark the positions of the conserved residues that are marked on the domain organization figures in the entries in Section II. The sequences of the following proteins are from the Swissprot database unless otherwise indicated and the database accession number and residue numbers are given in brackets. CD4d2, human CD4 precursor domain 2 (P01730, 123-204); NCAM-L1d3, mouse neural cell adhesion molecule L1 precursor domain 3 (P11627, 243-331); MAGd4, rat myelin associated glycoprotein precursor domain 4 (P07722, 327-412); Amalgam d3, Drosophila amalgam protein precursor domain 3 (P15364, 231-327); NCAM, chicken neural cell adhesion molecule precursor (P13590, 203-295); CEA d4, human carcinoembryonic antigen precursor domain 4 (P06731, 321-410); IgFcRII, mouse IgG FcRII precursor domain 1 (P08101, 37-116); CD2d2, human CD2 precursor domain 2 (P06729, 127-203); CD3 epsilon human CD3e precursor (P07766, 29-117).

domains, β 2-microglobulin ¹², MHC Class I antigen α 3 domain ³⁴, CD4 domains and 2 ^{8,9} and CD8 α ³⁵. The structure of domain 1 of CD2 has recently been determined by NMR ¹⁰. These structures show that the IgSF domains characterize by sequence similarities over about 100 amino acids correspond to structural units with distinct folding patterns referred to as the Ig-fold [reviewed in ref. 12]. The Ig fold consists of a sandwich of two β sheets, each consisting of antiparallel β strands of 5-10 amino acids with a conserved disulphide between the two sheets in most but not all domains. The sequence similarities are mainly found at the positions in-pointing residues in the β strands with considerable differences in the loops that connect the strands and the out-pointing residues on the faces of the β sheets. The core of the fold is made up of three β strands labelled ABE and GFC and the positioning of these is shown in the various folds illustrated in Fig. 9. The folds vary considerably in length in the middle of the sequence with Ig V-domain fold being the archetype for the longer fold. The extra sequence in comparison with domains forms an additional pair of β strands (C' and C'' in Fig. 9) and the connection between these forms the second complementarity determining region in antibody and TcR V-domains.

Figure 13. Integrin α chains. Residues identical in 3 out of 4 of the sequences are boxed. The sequences of the following proteins are from the Swissprot database and the database accession number and residue numbers are given in brackets. CD49a, rat integrin $\alpha 1$ precursor (P18614, 25–1172); CD49b, human integrin $\alpha 2$ precursor (P17301, 26–1161); CD51, vitronectin receptor integrin αV precursor (P06756, 27–1023); CD49d, integrin $\alpha 4$ precursor (P13612, 36–1013). The extracellular and transmembrane regions are shown.

CD49A FCVSPNVVVKNSMSFGPVEDMFGYTVQQYENBEQ...KVVLLQSPVLVQPKA...RTGVYKCPVGRERAMPVVKLDLPV
 CD49B CLAYLVNGLPEAKIFSGPSSBGFYAVQQGIPNPKG...NWLVLGSPVSWGPPEN...RMGDVYKCPVDLSTA-TCEKLNQOT
 CD51 LCRAFNLDDVDSPAEYSGPEGSYFCVAVDFVSSASRRMFLVYGAPKAN-TTPQGIYEGGQVTKC...DWSSTRRCQPTLEDA
 CD49D TORPYNVDTESALLYQGPHNTLPGYSV...VLHSHGANKRWLLVGAAPTANWLANASVINPGATVRCRIGKNPGQTCEQLQDGS
 NTSIP...NVTEIKENMTFGSTLVYNP...NGGFLACOPLYAYRCGHLY...TTGVSDVSPVTFQVVMSPAP-VQEGSTQ-
 STSIP...NVTEIKENMTFGSTLVYNP...NGGFLACOPLYAYRCGHLY...TTGVSDVSPVTFQVVMSPAP-VQEGSTQ-
 TGNRYAKDDPFLKSNHSLGLILTRNMG-TGGFLTCGPLWAQCCGNQY...TTGVSDVSPVTFQVVMSPAP-VQEGSTQ-
 PNOEPCGK-TCLERDNQWLGVTSRQPGENGIVTCG...RSKQDKILACAPLYHWKTEM...KQEREPTVTCFLQDGTGTVY...APCRSQD
 LDIVVLDDGNS-195 Amino acids...SGTGFSAHYSSQNDILMLGAVGAFQWSTIVQKTSHGHLI...TPHNTTFTTEPAKMNE
 IDVVVVCDEN-195 Amino acids...SQVGFSAHYSSQNDILMLGAVGAFQWSTIVQKTSHGHLI...TPHNTTFTTEPAKMNE
 IDADGGGF...CQGGFSDFT-KADRVLLGPGSFYVWQQQLISDQVAEIVSKYDPNVYSIKY...FPKQAFDQILQDRNH-
 VKKGFENFA...CQAGTSSFT-KDLIVMGAPGSSSYWTCGLF...VYNJTNNKKYKAFLDKQNV
 PL-ASYLGTVNSATIPGD...VLYIAGQPRYNNHT-GVQVLYKMGEDQMNIQTLCGEGISYFGSVLTITIDIDKSYTDLVLLVGA
 ...SSYLGYSAVAAISTGES...THFYAGAPRANYTGIVLYSVNENQNTIVQAHGDDQIGSYFGSVLTITIDIDKSYTDLVLLVGA
 IFDSSYLGYSAVAGDFNGDIDDFVSGVPAARLGMVYIY...DGKNNMSSLYNFTGEOMAAVYFGFVAAATDINGDDYAD-VFVGA
 KP-GSYLGYSVAGHFRSQHTTEVVGAPQHEQ-IGKAYLPSIDEKELNHEMKKKLGSYFGSVLTITIDIDKSYTDLVLLVGA
 PMYM...GTEKEEKGKVVVYAVNQTRFEYQMSLEPIRQTCSSSLKDNSCTKKNKNEPCQARFGTATAAVKDLNVDFGNDVVLGA
 PMYM...GTEKEEKGKVVVYAVNQTRFEYQMSLEPIRQTCSSSLKDNSCTKKNKNEPCQARFGTATAAVKDLNVDFGNDVVLGA
 PLFM...DRGSDGKLEVGQVQVSVLSLQASGDF...EGEGJENTRFGSIAIA...SDINMDGNDVVLGA
 PM...QSTIREGRVYFVYINSGGAV...MNAMEINLVGSDXYARFGESLIVNLGIDINDGDFEDVVLGA
 PL-BDDHAGAVYIYNGSOKTIRBAYAQRPSSGDDO...KTLKFFGQSIIHDEMNLNGDGLTIDVTIGGLG...GAALFWARDAVAVKVT
 PL-BENQNSGAVYIYNGSOKTIRBAYAQRPSSGDDO...KTLKFFGQSIIHDEMNLNGDGLTIDVTIGGLG...GAALFWARDAVAVKVT
 PYGDEDKAGIYIENGRTQTLNAPVSLLEQWAA-RSMPPSPFGYSMKGATIDKNGYVPLIVGAFGVDRALYLRARPVITVNAO
 P-DEDDLQGAITYLNGRADQISSTFSQRTEDLQIS-KSLSM-FGQSISGQIDADNNGYVDAVAGAFRSDSAVLLRTRPVVIVDAS
 MNFBPNKKNVQKXNCRVEGK...ETVCINATMCFHVKLKSKEDSIYEADLQVRYVTL...SLRQISRSFSGTQERKIQRNITVR
 ASFTPEKTLNKNNAQITLK...LCEASAKRPTKQN-NQVAIVNITDADGSSRSTVSGGLHKKENNERCLOKNNVVR
 LEVYPSILNQDNQNTSLPGTALKVSFPNVRFCLEKAGRGV...LPRKLNFGVLLDLKLKQKGAIRRALFLYSRSPSHSKNMTI...
 LS-IPESVNRRTKFDG...VENGWPSVQIDITLCEFSYKCKEY...PGYIVLVNMSLD-VNRKKAESPFRFYSNNGTSVDVITGSIQV

ESB...CIRHSFYMLDKKHDFQD...SVRVTLDFNLTD...PENGPIVLDALPNSVHEHIPFAKDCGKNKERCIISDL
 QAQS...CPEHIIYIQEPSDVN...SLDLRVDISLEN...POTSPALLEAYSETAKVESIPPHKDCGEGDLGICISDL
 SRGLMQCEELIAYLRDESEFRKLTPTITIMEYRLD...YRTAADTGLQILNQFTPANISRAHILLCCGEGDNVCKPKL
 SSREA-NCRTIQAQFMK...DVRDILTPIQIEAAYHGLPHVISKRSTEEBPLQPIQLQKKEKDKIMKKTINFARFCAHENCISDL
 TLMV...STTEKSLTVKSHQDKENVSILTVKKNKGDSEYNTVQHSNPLIFSGI...BGIQKDSCEBQ...NITCRVVG
 VLDD...RQIPAAQEQPHIVSNQNKRLTPSVTLKKNKRESAYNIGIVDFSENLFFHASF-SLPVDTGTEVTCQVAASQK...SVACDVG
 EVSV...DSQKKIYIGDNDNPLTIVKAK...NQEGAGYAEALIVSIPQLADDFIGVVRNNEALARLSCAFKNTRQVVCDLG
 QVSAKIGFLKPHENKTYLAVGSMKTLMNLVSLFNAGDDAYBTTLNVKLPVQLYFIKILELE...KQINCEVTDNSGVVQLDCSIG

ESE... CIRHSFYMLDKHDFQD... SVKVTLDNFNLT... PENGPVLDDALPNSVHEHIPPFAKDCGNKERCISIDL
 QAQS... CPEHIIYIQEPSDVVN... SLDRVDISLEN... PGTSPALAEAYSETAKVFSIPFHKDCEDOLCISIDL
 -SRGOLMQCEELIAYLRDESEFRDKLTPTITIMBYRLD... YRTAADTGLQPIINQFTPANISRQAIIILLDCEDNVCKPKXL
 SSREA... NERTHQAFMRK... DVVRDILTPIQIEAAYHGLGPHVISKRSSTEEFPLQPIILQQKKEKDIMKKTINFARFCAHENCADL
 TLNV... STTEKSLITVKSQHDKEFNSVLTVMKQDSJAYNTRTVMQHSNMLIFSGI... BEIQKDSQESNQ... NITCRVVG
 EVSV... RQIPAAQEPFIIIVSNQNKRRITFSTTLKNNKESAYNIGIVVDSENFASF... SLPVDGTEVTCQVAAASQK... SVAACDVVG
 FVSV... DSDQKKIYIGDDNPLTLIVKQA... NQGEQAYEAEIIVSIPLQAFIGVVRNMALARLSCAFKTENQTRQVVCDDLIG
 QVSAKIQFLKPHENKTYLAAGSMKTLMLNVSLFNAGDDDAYETTLHVVKLPVGLYFIKILELEB... KQINCCEVTDNSGVVQLDCSIG

YPFLRAGETVTFKIIIFQFN... TSHLSENAIIHLSATSDS... EEPLESNDNEVINISITPVKYEVLQFYSSASAEHHSVAAANETIPEP
 YPALKRREQQVTFITINFDNF... LQNLQNASLSFQALSESQEB... NKADNLVNLKLELLYDAEIHILTRSTNINFYEISSDONVPSI
 NP... MKAGTQLLAGLRISVILQSEMDTSVKFDLQIQSNIFDKVSPVVSXKVDLAYLAVEIR... GVSPPDHIFLPIPNWEHKN
 YIYVDHLSRIDISFLDVSSLSRAEBEDLSITVHATCENEEB... MDNLKHSRVTVALPLKYEYK... LTVHIOFVNPTSFVYGSN
 INSTEDIGN... EINVPYTIKRRGHPPMPPELQLSISFFN... LTADGYPVLYPIGWS... SSDNVNCRPRSLEDPIFGTINSGKKMTISKSE
 VHSFEDVGP... KFIPSLKVTIGSVPVSMATVILHIPO... YTKBNP... LMYLTGVQ... TDKAGDISCNADINPLKIQTSSSVSSEKSE
 PETEEDVGPV... VQHIYELRNNGPSSFSKAMHLQWPKYNNNT... LLYILHYD... IDGPMNCTSDMEINPLRIKISSLQTTKND
 DENEPETCMVEKMNLTFFVINTGNMAMPNVSVETMVPNSPSPQTDKLFNILDVQTTTBOCHFENYQRYVCALIEQQKSAMQTLKQIV
 VLKROTIQD... CASS... TCGVATITCSLLPSDLSQ... NVISLLWKPTFIRAHFSSNLTLRQ...
 NFRHTKELN... CRTASCS... NVTCWLKDVHMKGEYFVNVTTRIMNOTFASSTFQVQLTAA...
 TVAGQGERDHLITKRDALSEGDIHTLQCGVACQL... KIVCQVGRLDRGKSAILEVYKSLWTEIFMKNENQNSYSLKSSASPNV
 RFLSKTDKRLLYCIK... ADPIICL... NPLCNFGKMESSOKEASHIQLEGRPSILEMDETSALKFEIRATG...
 BLKSENSSSLTSSSNRKRRELAIQISKDGLPGRVPLWVILL... SAFAGLLLLMLILLALWXLIGFFKR
 RINTYNPRIYVIEDN... TITIPLMIMKPDKAEVPTGVIIIG... SIAGTILLALVALLWXLIGFFKR
 IEPYKKNLPIEDITNSTLVTTNTVWTGIIQAPMPVPVWVILL... AVLAGTILLALVALVFMVYRMGFFKR
 ...FPEPNPRVIELNKNDENV... AIIVLEGLHHQRPKRYFTIIVISSILLGLITVLLISYVMWKAIGFFKR

cytoplasmic domain

In the IgSF there are limited sequence patterns in β strands B, C, E, and F that are common across the superfamily (Figs 10-12) and other limited patterns that allow a subdivision of the domains. Ig and TcR V-domains have a characteristic pattern in the region leading into β strand F of Asp-X-Gly/Ala-X-Tyr-X-Cys. The receptor C-domains have a characteristic pattern between β strands B and C of GlyPheTyrPro and another on the COOH-terminal side of β strand F of Cys-X-Val-X-His. The Ig, TcR and MHC antigen C-type domains all share the same types of sequence patterns and are referred to as the C1 set within the IgSF. With the sequencing of various cell surface molecules a third category of domains became evident, namely domains of length similar to C-domains but with some of the sequence patterns of V-domains. These domains are referred to as the C2 set². They have V-type patterns in the β strand E to F region, a pattern of Pro-X-Pro is relatively common between β strands B and C and the pattern Cys-X-Ala-X-Asn is common after β strand F. CD4 domain 2 is a C2-set sequence and its structure is classified in terms of sheet assignments labelled as ABE/GFCC'. This is in comparison to ABED/GFC for C1-set sequences and ABED/GFCC'C" for V-set sequences. That is, for C2-set sequences the middle β strand may be generally in line with the GFC β sheet rather than the ABE sheet as is the case with antibody C-domains. The points about conserved patterns and the positioning in β sheets are made evident by comparing the sequence alignments in Figs 10-12 with the folding patterns for the domains in Fig. 9. The sequence alignments are discussed in more detail in ref. 2.

Integrin superfamily (Figs 13 and 14)

The integrins are a large family of related proteins that all share a heterodimeric structure with α and β chains that both traverse the lipid bilayer. There are at least 20 α and eight β chains which can be found in various but not all possible combinations. Sequence similarities are seen within the α and β chain across all the integrin types (Figs 13 and 14). The integrins are known to be involved in cell interactions and include receptors for the extracellular matrix proteins fibronectin and vitronectin and for cell surface molecules ICAM-1 and ICAM-2. The integrins have been extensively reviewed elsewhere including a companion volume in this FactsBook series³⁶. They are expressed on many different cell types; the CD11/CD18 family and the CD49 very late activation antigen family (VLA) are expressed mainly on leucocytes. The sequence similarities in this family are described in more detail in refs 37-39.

This family of related proteins does not contain other domain types apart from the $\beta 4$ integrin that contains two FN type III SF domains in the cytoplasmic region (see ref. 33 and section on fibronectin type III SF domains).

Figure 14. (opposite) Integrin β chains. Residues identical in 3 out of 3 of the sequences are boxed. The sequences of the following proteins are from the Swissprot database and the database accession number and residue numbers are given in brackets. Beta 1, human fibronectin receptor, integrin β_1 precursor (P05556, 26-752); Beta 2, human integrin β_2 (CD18) precursor (P05107, 024-724); Beta 3, human integrin β_3 (CD61) precursor (P05106, 30-742). The extracellular and transmembrane regions are shown.

Figure 14. Integrin β chains

Beta 1	RGLKANAKSGQBCIQAGPNCQWCTNSTF.LQEGMPTSA	RCEDLEA	KKKGPPDD	ENPRGSKDIKKKNVT
Beta 2	ECTKFKVSSCBEDIESGPGCTWCQKLN.F.TOPGPD	IRCDTRPQ	LMRGCADD	IMDPTSLAETQEDHNG
Beta 3	ICTTRGVSSQCCGLAVSPMCAWCSDEALPLGSP...	RCDLKEN	LKDNCAPE	IEFVSEARVLEDRPLS
	NRSGTAEKLKPEDIHQIQPQQLVRLRLRSGEPQFTLKFKRAED	YPI	DIYYLMDLSYSMKDD	LENVKS LGTDQWNE
	QK.....QLSPQKVTLRLPQQAANVTFRRAKGYPI	DIYYLMDLSYSMKDD	LENVKS LGTDQWNE	LRNVKLLGGDLRA
	DKSODSS.....QVTVSPQRIALRLPDDSKNFISQVRQVED	YPI	DIYYLMDLSYSMKDD	LENVKS LGTDQWNE
	MRRITSDERLGFGEFVENITMPYISTIP.AKLRNPQTS.EQNTTTP	SYKKNV	ITNKGEVFNE	LVGKQR IISGNLD
				ITNKGEVFNE LVGKQR IISGNLD

that are
t allow a
attern in
eptor C-
eTyrPro
. The Ig,
sequence
ncing of
, namely
tterns of
e V-type
common
1 after β
in terms
ED/GFC
r C2-set
et rather
ts about
mparing
mains in

odimeric
e at least
possible
cross all
ed in cell
onectin
integrins
e in this
pes; the
VLA) are
mily are

part from
ic region

3 of the
rom the
ibers are
recursor
24-724);
acellular

Figure 14. Integrin β chains

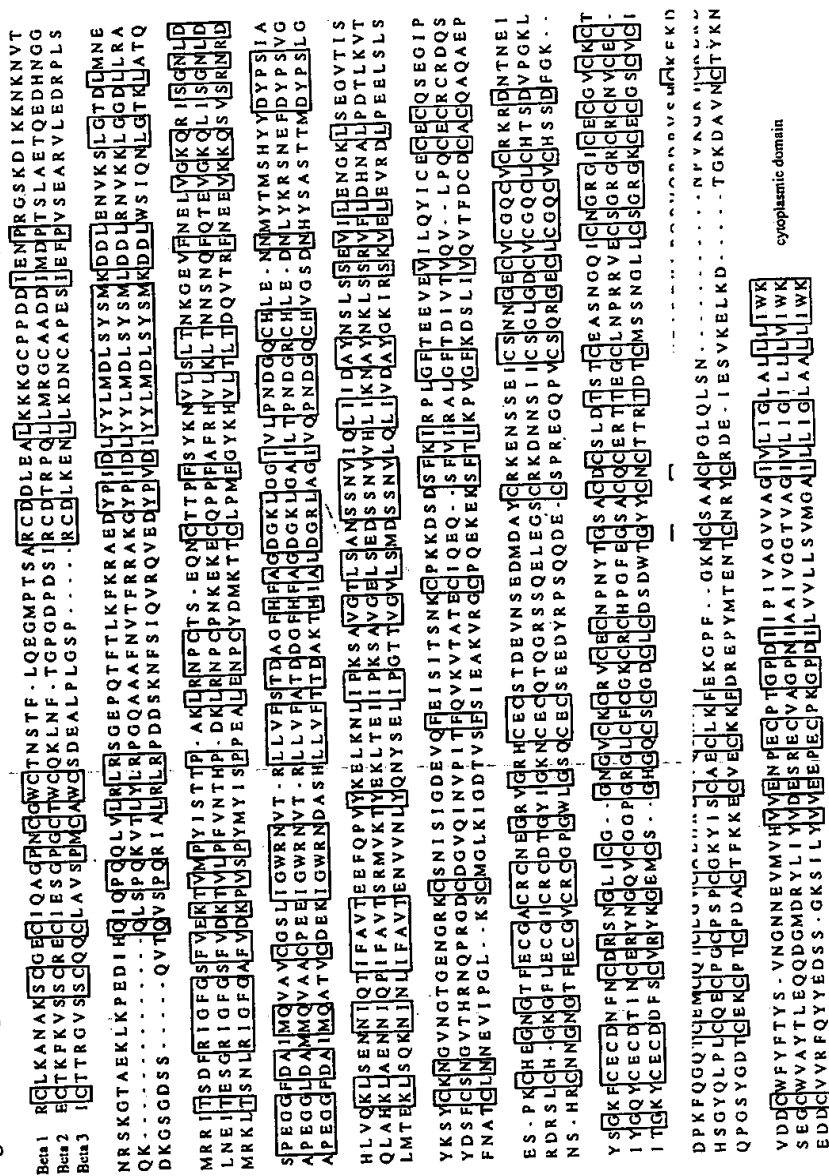
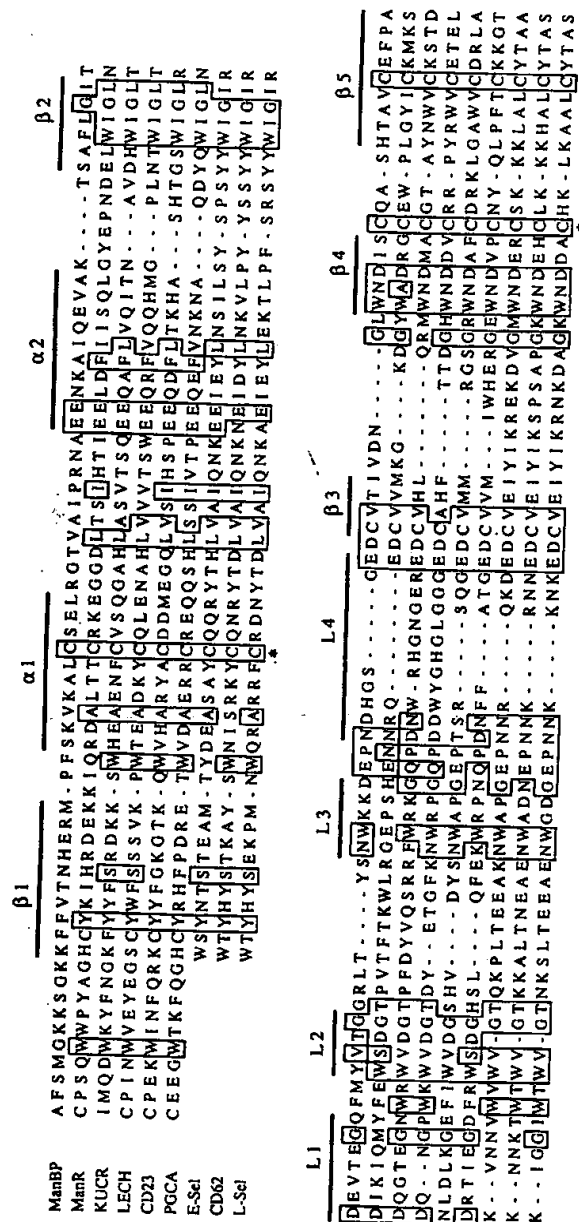
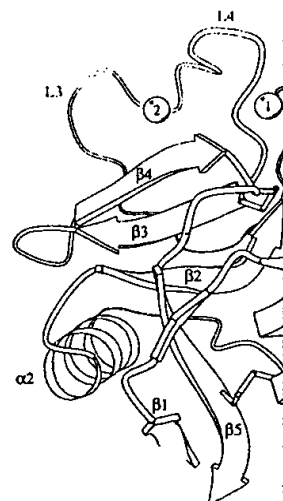


Figure 15. Lectin C-type superfamily domains. Residues identical in 4 out of 7 of the sequences are boxed. The positions of the β strands (β), α helices (α) and loops (L) determined for the structure of the rat mannose binding protein 49 are shown above the sequences. The asterisks mark the positions of the conserved residues that are marked on the domain organization figures in the entries in Section II. The sequences of the following proteins are from the Swissprot database and the database accession number and residue numbers are given in brackets. ManBP, rat mannose binding protein A precursor (P19999, 117-238); ManR, human mannose receptor precursor (P2897, 362-490); KUCR, Kupfer cell carbohydrate binding receptor (P10716; 412-540); LECH, rat hepatic lectin-1 or asialoglycoprotein receptor-1 (P02706; 152-279); CD23, low affinity IgE receptor (P06734; 163-286); PGCA, cartilage specific proteoglycan core (P07897; 1914-2038); E-Sel, E-selectin or ELAM-1 precursor (P16581; 22-142); CD62 or granule membrane protein 140 or P-selectin precursor (P16109; 42-162); L-Sel, human leucocyte adhesion molecule or L-selectin precursor (P14151; 39-159).



Lectin C-type superfamily (P) This family of lectin domain shown to bind carbohydrate found in a number of lectin hepatocyte galactose receptor binding proteins in two in Lectin C-typeSF domains are bind carbohydrate, such as leucocyte adhesion molecule as L-selectin and the low carbohydrate binding for the cartilage proteoglycan core p

Two groups of lectin C-type lectin C-typeSF domain plus completely within one exon the selectins, the lectin domain include the COOH-terminus the genetic region encoding intron boundary and thus if might function to produce terminus. If the insertion membrane protein the lectin little functional relevance. would form a new extracellular properties. The recent cDNA the first example of a protein There is no information as yet



Lectin C-type superfamily (Figs 15 and 16)

This family of lectin domains are termed C-type because some members have been shown to bind carbohydrate in a Ca^{2+} dependent reaction ⁴⁰. This domain has been found in a number of lectins such as the Kupffer cell fucose/galactose receptor, hepatocyte galactose receptor, mannose binding protein from plasma, and galactose binding proteins in two invertebrate species, the flesh fly and sea urchin ⁴¹⁻⁴³. Lectin C-typeSF domains are found in a number of proteins not originally known to bind carbohydrate, such as the proteoglycan core protein ⁴⁴, an endothelial leucocyte adhesion molecule (E-selectin), and leucocyte cell surface antigens such as L-selectin and the low affinity Fc receptor for IgE (CD23). In some cases carbohydrate binding for the lectin C-typeSF domain has been established, e.g. cartilage proteoglycan core protein ⁴⁵.

Two groups of lectin C-type domains can be distinguished. The L-selectin has the lectin C-typeSF domain plus about 10 residues of the signal sequence contained completely within one exon with phase 1 intron boundaries ⁴⁶. In cases other than the selectins, the lectin domain is usually found spread over three exons which also include the COOH-terminus of the protein and the 3' untranslated sequence ⁴⁷. In the genetic region encoding the NH_2 -terminal side of the exon there is a phase 1 intron boundary and thus if this exon was inserted into an intron of another gene it might function to produce a new protein with the lectin domain at the COOH-terminus. If the insertion occurred after a hydrophobic region in a Type I membrane protein the lectin exon would be cytoplasmic, and thus presumably of little functional relevance. Conversely if it were inserted in a Type II protein it would form a new extracellular COOH-terminal region with carbohydrate binding properties. The recent cDNA sequence of the macrophage mannose receptor ⁴⁸ is the first example of a protein containing multiple lectin repeats, with eight in all. There is no information as yet on its genomic structure.

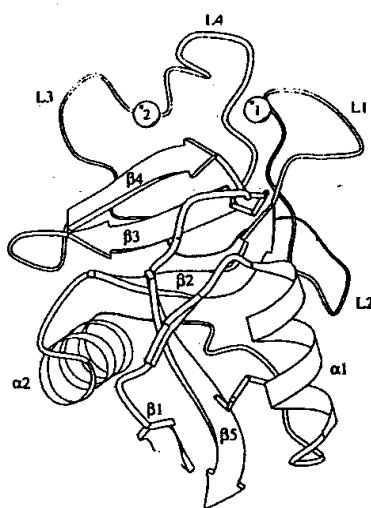


Figure 16. The folding pattern of a lectin C-typeSF domain. Ribbon diagram of the lectin C-typeSF domain from the rat mannose binding protein ⁴⁹. The β strands are shown as broad arrows pointing from the amino to carboxy direction, α helices as coiled ribbons and the connecting loops as thinner lines. The labelling of the β strands ($\beta 1-5$), α helices ($\alpha 1-2$) and loops (L1-4) corresponds to that in the sequence alignments in Fig. 15. The numbers 1 and 2 refer to the position of the two holmium ions that are known to stabilize this region that contains a high proportion of nonregular secondary structure ⁴⁹.

66

As well as the different patterns that differ between the two, it is evident in the sequence that the amino acid residue at the NH₂-terminus of the lectin, whilst in the lectin terminus that shows a

The structure of a has recently been determined (Fig. 16). The structure contains non-regular crystal structure. The domains for the other

Lectin S-type superficial
Galactoside binding protein
to contain a sequence
been termed S-type
groups⁴⁰. These are fr
strong sequence sim
However, analysis of
requirement for a red
groups⁵⁰.

CD42a
CD42b
A2g-1
A2g-2
A2g-3
PG-1
PG-2
PG-3
Chao-1
Chao-2
CYCL-1
CYCL-2

Figure 18. *Leucine-rich sequences are boxed, that are marked on the sequences of the otherwise indicated as given in brackets. CD56-79; CD42b, human A2g-1, A2g-2 and A2g-134-157, 158-181 and proteoglycan II precursor Chao-1 and Chao-2, I (respectively), CYCL-1 891-913 respectively).*

As well as the differences in intron/exon organization there are patterns that differ between these two groups of lectin C-typeSF domains evident in the sequence alignments shown in Fig. 15. There is a characteristic residue at the NH₂-terminus of the selectins E-selectin, CD62 (P-selectin), whilst in the other group there is a longer patch of sequence at the terminus that shows a conserved sequence pattern of Cys-X-X-X-Trp.

The structure of a lectin C-typeSF domain from a rat mannose binding has recently been determined by X-ray crystallography⁴⁹ and the fold is shown in Fig. 16. The structure is unusual in that it contains two regions, one of which contains non-regular secondary structure stabilized by two holmium ions, and the other contains both β sheet and α helix; this is unusual for a cell surface molecule. The other cell surface molecules consist solely of β structure.

Lectin S-type superfamily (Fig. 17)

Lectin S-type superfamily (Fig. 17) Galactoside binding proteins have been sequenced from several species and shown to contain a sequence pattern different from that of the lectin C-type superfamily. These have been termed S-type because the first examples contained free accessible thiol groups⁴⁰. These are found both intracellularly and extracellularly and a strong sequence similarity is found in the Mac-2 leucocyte antigen (Fig. 17). However, analysis of protein produced by recombinant DNA technique shows no requirement for a reducing environment for lectin activity and no accessible thiol groups⁵⁰.

[illegible]

Figure 18. Leucine-rich glycoprotein repeats. Residues identical in six sequences are boxed. The asterisks mark the positions of the conserved residues that are marked on the domain organization figures in the entries in Table II. The sequences of the following proteins are from the Swissprot database unless otherwise indicated and the database accession number and residue number are given in brackets. CD42a, human platelet glycoprotein IX precursor (P14770, 56-79); CD42b, human platelet glycoprotein 1B β chain precursor (P1224, 60-83); A2g-1, A2g-2 and A2g-3, human leucine-rich $\alpha 2$ glycoprotein (LRG) (P12750, 134-157, 158-181 and 182-205 respectively); PG-1, PG-2 and PG-3, human bone proteoglycan II precursor (P07585, 86-109, 110-133 and 134-157 respectively); Chao-1 and Chao-2, Drosophila chaoptin precursor (P12024, 132-155 and 156-179 respectively); CYCL-1 and CYCL-2, yeast adenylate cyclase (P08678, 58-890 and 891-913 respectively).

Leucine-rich repeats (LRR) or leucine-rich glycoprotein (LRG) repeats (Fig. 18)

The leucine-rich repeat (LRR) is characterized by a pattern of conserved residues including about 5 or 6 leucines and some other residues in a tightly defined repeat of 24 residues (Fig. 18). It is found both intracellularly and extracellularly in a variety of species including *Drosophila* and yeast and has also been found in the platelet glycoproteins CD42a and CD42b. It often occurs in an array of tandem repeats. For instance there are nine repeats in the leucine-rich glycoprotein where the repeat was first noted, 26 repeats in the yeast adenylyl cyclase, 10 in the proteoglycan protein⁵¹ and three in the trkB protein⁵². In some cases some sequence similarity is observed beyond the alignments shown and an alternative alignment may start from the conserved Pro position which is in the centre of the alignment shown. The α chain of the CD42b contains seven LRRs and these, together with all the remaining coding sequence, are encoded by a single exon⁵³. Thus, this repeat is not generally coded by single exons in this case. The LRRs have been found in diverse proteins and they have been implicated in the specificity of hormone binding to gonadotropin receptors⁵⁴ and in the interaction between yeast adenylyl cyclase and RAS proteins⁵⁵.

The leucine and other residues in LRRs form an amphipathic sequence which could be involved in protein-protein or protein-lipid interactions. One of the repeats of the *Drosophila* chaoptin protein has been synthesized. This peptide is soluble in aqueous solution but will bind to phospholipid vesicles where it forms predominantly a β structure. It has been suggested that protein segments containing tandem repeats may also form amphipathic β sheets⁵⁶.

Link superfamily (Fig. 19)

Two link superfamily domains were originally noted in the link protein that binds hyaluronic acid⁵⁷. This protein also has one IgSF domain. Subsequently, a further four linkSF domains were observed in the proteoglycan core protein that has a chondroitin sulphate binding site. This protein also contains one IgSF domain, a CCPSF domain and a lectin C-typeSF domain⁴⁴. There is also a single linkSF domain in the CD44 antigen which is known to bind to hyaluronate^{58,59}.

Low density lipoprotein receptor (LDLR) superfamily (Fig. 20)

The LDL receptor contains seven domains of about 40 amino acids with six conserved cysteine residues that have been called LDLRSF domains⁶⁰. The LDLR also contains three EGFSF domains. LDLRSF domains have also been found in other proteins, notably some complement components such as C6, C9 and factor I. In the LDLR, four of the LDLRSF domains are each encoded by one exon whilst the other three are encoded by a single exon⁶¹. In the LDLR mutational analysis has indicated that the LDLRSF domains are important in the binding of some lipoproteins but otherwise the function of this domain type is not known⁶². The structure of the LDLSF domain has not been determined to date.

Ly-6 superfamily (Fig. 21)

The Ly-6 antigens are a group of leucocyte antigens first identified in the mouse that consist of 70–80 amino acids containing 10 Cys residues^{63,64}. Southern blot analysis indicates that many Ly-6-related genes are present in the mouse and of these, 10 distinct genes have been identified⁶⁵. The Ly-6 antigens are expressed in non-lymphoid tissues, for example, kidney, as well as on leucocytes. Homologues

Figure 19. Link superfamily domains. Residues identical in four or more sequences are boxed. The asterisks mark the positions of the conserved residues that are marked on the domain organization figures in the entries in Section II. The sequences of the following proteins are from the Swissprot database and the database accession number and residue numbers are given in brackets. CD44, human CD44 antigen precursor (P16070, 32–123); CORE1, CORE2, CORE3 and CORE4, rat cartilage-specific proteoglycan core protein precursor (P07897, 152–251, 253–353, 486–585 and 587–687 respectively); LINK1 and LINK2, rat proteoglycan link protein (P03994, 143–242 and 244–339 respectively).

CD44
G V F H - V E K N G R Y S I S R T E A A D L C K A F N S T L P T M A Q M E K A L S I - G F E T C R Y G L F
A T U A P I V I R A A S T I R Y T I D E D R A Q R A C L Q N S A I I A T P E Q L Q A Y E D - G F H Q C D A G W
G M I D M I C I S A G W

Figure 19. Link superfamily domains. Residues identical in four or more sequences are boxed. The asterisks mark the positions of the conserved residues that are marked on the domain organization figures in the entries in Section II. The sequences of the following proteins are from the Swissprot database and the database accession number and residue numbers are given in brackets. CD44, human CD44 antigen precursor (P16070, 32-123); CORE1, CORE2, CORE3 and CORE4, rat cartilage-specific proteoglycan core protein precursor (P07897, 152-251, 253-353, 486-585 and 587-687 respectively); LINK1 and LINK2, rat proteoglycan link protein (P03994, 143-242 and 244-339 respectively).

GLDLR-d1	GTAVG--DR	CERNEP	QCQDG	--KCI	SYKWV	CDGSA	ECQDGSD	ES
GLDLR-d2	ETCLDL	--VTCK	SDFSC	GGVR	NRCIP	QFWRC	DGQV	CDN
GLDLR-d3	--GCLP	--KTCS	QDFRCH	DG--KCI	SRQFV	CDSD	DRD	LDGSD
Comp 9	EQALP--	SECS	SIEFTCES	G--ACI	KRLSC	ZN	CDYD	CEGSD
Hemo. Linker	DELEG--	NGCE	PRHFQ	CGGSA	MECIS	DLT	CDGSP	CANGSD
Factor 1	ELCLC--	KACQ	GKFRHCK	SG--VCI	PSQYQ	CNG	EVDC	ITGDE
Comp 7	RGCPT	E--EGC	G--GERF	CFSG--QCI	SKSLV	CNG	DCD	DESD
Comp 6	LLCKIE	BA	AD	C--K	NKFRC	DSG--R	CLARK	LE
Comp 6								

Human CD59
Mouse Ly-6A
Mouse Ly-6C
UPAR-1
UPAR-2
UPAR-3
Squid Sgp2

LQCYNC	PN	-PTAD	CKT	--AVNC	SSD	FDA	C	--LITK	A	GLQ
LECYQC	YQ	VVPF	ETS	CP	S	PDP	DGVC	--VTQE	A	AYVI
LCYC	EYV	PIET	SCPA	--VTCRA	SDGFC	--TAQN	L	TEB	I	
LRCMQ	CKT	-NGD	CRV	--EECAL	GQD	LCRT	T	LVRL	WEEG	
LECI	SCGS	-DMS	CERGR	--HSQL	CRSP	EEQC	LD	VVTH	WIJ	EQE
RQCY	SKGN	-SHGC	SSBE	--TFLLD	CRCP	PMNQ	C	LVAT	GTHE	
IKCFV	CNSY	-HQDC	GDWF	NATHSV	HQCEPS	QRKR	K	LVQQ	I	KLD
V	--YNKC	WK	--FBHC	NFNDV	TTTLR	KREN	-ELTY	YCC	KKX	DLCLCN
VVD	SQTR	KVKNN	KL	--PICPP	-NIESMEI	LOTKVNV	KTS	CCQ	EEB	EDLCN
ED	SQRR	KKLT	RQCL	--SFCA	OVP	PKDPNIR	--ERT	SCC	EEB	EDLCN
EE	--LELVE	KSTH	--SEKTNR	TLTD	SYRTGL	KITS	TV	EVY	CG	DLCLCN
EE	GRP	KDDRHL	RQCGY	--LPG	CPG	SNGFHN	DDTF--	HF	LKCC	TKSGCN
EE	NQSMYVR	GGCAT	--SMQC	HAHL	ODAF	FSMNH	--IDVS	CCC	TKSGCN	
PEK	--WOVRY	--IROCA	EGGEIG	AYDOR	VCKDR	RGTSGVK	--MTY	CHC	QT	TEGCN

Figure 21. (opposite) Ly-6A more sequences are boxed residues that are marked in Section II. The sequence database unless otherwise noted. Mouse Ly-6A, (UPAR-2 and UPAR-3, I S12376, 23–99, 115–199) 2 residues 1–92⁶⁷.

Another member of receptor. This molecule and is also attached superfamily are shown remains to be determined domains of any other known for this superfamily.

The MHC superfamily
The MHC antigens
their membrane prox
terminal segments, i
and the $\alpha 1$ and $\beta 1$
similarity to IgSF s
independent structu
show weak sequenc

Figure 20. (opposite) Low density lipoprotein receptor (LDLR) superfamily domains. Residues identical in four or more sequences are boxed. The asterisks mark the positions of the conserved residues that are marked on the domain organization figures in the entries in Section II. The sequences of the following proteins are from the Swissprot database and the database accession numbers are given in brackets. LDLR, human low density lipoprotein receptor precursor (P01130, d1 20-59, d2 61-102, d3 103-141); Coelomocyte, rainbow trout complement C9 (P06682, 72-112); Hemo.Linker, marine worm extracellular haemoglobin linker 2 chain (P18208, 61-102); Factor D, human complement factor I precursor (P05156, 253-291); Comp 7, human complement C7 precursor (P10643, 77-116); Comp 6, human complement C6 precursor (P13671, 131-171).

Figure 21. (opposite) Ly-6 superfamily domains. Residues identical in three or more sequences are boxed. The asterisks mark the positions of the conserved residues that are marked on the domain organization figures in the entries in Section II. The sequences of the following proteins are from the Swissprot database unless otherwise indicated and the database accession number and residue numbers are given in brackets. CD59, human CD59 antigen (P13987, 26-95); Mouse Ly-6A, (P05533, 27-105); Mouse Ly-6C, (P09568, 2-102); UPAR-2 and UPAR-3, human urokinase plasminogen activator receptor (PIR; S12376, 23-99, 115-199 and 214-294 respectively); Squid Sgp2, squid glycoprotein 2 residues 1-92⁶⁷.

of the Ly-6 antigens have been found in the rat but not yet in humans the CD59 antigen has been identified as a member of the Ly-6 superfamily but seems too different in sequence from the mouse Ly-6 antigen to be a Ly-6 homologue. CD59 is a downregulatory control protein for human complement and also shows adhesion reactivity with the CD2 antigen⁶⁶. An invariant member of the Ly-6 superfamily has been isolated from squid optic tectal tissue^{64,67}. All the above molecules consist of a single Ly-6SF domain attached to the cell surface by a GPI anchor.

Another member of the Ly-6 superfamily is the urokinase plasminogen activator receptor. This molecule contains three domains separated by hinge regions and is also attached to the cell surface by a GPI anchor. The domain organization of this superfamily is shown in Fig. 21 and a tertiary structure for this domain type remains to be determined. No Ly-6SF domain has been found in combination with other exon structures known for this superfamily are not suited to exon shuffling (Table 1).

The MHC superfamily (Figs 22-24)

The MHC antigens and related molecules are members of the IgSF C1 set. However, their NH₂-terminal segments, including the $\alpha 1$ and $\alpha 2$ domains of MHC Class I heavy chain and the $\alpha 1$ and $\beta 1$ domains of the Class II α and β chains, show no sequence similarity to IgSF sequences⁶⁸ and the Class I domains are known to form an independent structural unit as shown in Fig. 22³⁴. The Class I $\alpha 1$ and $\alpha 2$ domains show weak sequence similarity to each other and form a similar fold containing a

family
the asterisks
domain
the following
n number and
protein
p 9, rainbow
1 giant
, human
complement
cursor

in three or
conserved
entries in
ssprot
umber and
1 (P13987,
102); UPAR-1,
eptor (PIR;
id glycoprotein

ther species. In
y-6 superfamily
ns to be a Ly-6
omplement and
tebrate member
central nervous
nain attached to

gSF by virtue of
ever, their NH₂-
ss I heavy chain
ow no sequence
own to form an
and $\alpha 2$ domains
fold containing a

Figure 20. (opposite) Low density lipoprotein receptor (LDLR) superfamily domains. Residues identical in four or more sequences are boxed. The asterisks mark the positions of the conserved residues that are marked on the domain organization figures in the entries in Section II. The sequences of proteins are from the Swissprot database and the database accession numbers are given in brackets. LDLR, human low density lipoprotein receptor precursor (P01130, d1 20-59, d2 61-102, d3 103-141); Ctrout complement C9 (P06682, 72-112); Hemo.Linker, marine worm extracellular haemoglobin linker 2 chain (P18208, 61-102); Factor complement factor I precursor (P05156, 253-291); Comp 7, human complement factor 7 precursor (P10643, 77-116); Comp 6, human complement C6 precursor (P13671, 131-171).

family
The asterisks
the domain
the following
ion number and
lipoprotein
np 9, rainbow
m giant
-I, human
complement
precursor

Figure 21. (opposite) Ly-6 superfamily domains. Residues identical in three or more sequences are boxed. The asterisks mark the positions of the conserved residues that are marked on the domain organization figures in the entries in Section II. The sequences of the following proteins are from the Swissprot database unless otherwise indicated and the database accession numbers are given in brackets. CD59, human CD59 antigen (P13987, 26-95); Mouse Ly-6A, (P05533, 27-105); Mouse Ly-6C, (P09568, 7-102); UPAR-1, UPAR-2 and UPAR-3, human urokinase plasminogen activator receptor (PIR; S12376, 23-99, 115-199 and 214-294 respectively); Squid Sgp2, squid glycoprotein 2 residues 1-92⁶⁷.

l in three or
conserved
e entries in
vissprot
number and
en (P13987,
7-102); UPAR-1,
ceptor (PIR;
uid glycoprotein

of the Ly-6 antigens have been found in the rat but not yet in humans the CD59 antigen has been identified as a member of the Ly-6 superfamily but seems too different in sequence from the mouse Ly-6 antigen to be a homologue. CD59 is a downregulatory control protein for human T cells and also shows adhesion reactivity with the CD2 antigen⁶⁶. An antigen of the Ly-6 superfamily has been isolated from squid optic tectum tissue^{64,67}. All the above molecules consist of a single Ly-6SF domain attached to the cell surface by a GPI anchor.

Another member of the Ly-6 superfamily is the urokinase plasminogen activator receptor. This molecule contains three domains separated by hinge regions and is also attached to the cell surface by a GPI anchor. The domain organization of this superfamily are shown in Fig. 21 and a tertiary structure for this domain type remains to be determined. No Ly-6SF domain has been found in the domain organization of any other superfamily and this may be because the domain organization known for this superfamily are not suited to exon shuffling (Table 4).

other species. In
Ly-6 superfamily
gens to be a Ly-6
complement and
vertebrate member
l central nervous
main attached to

The MHC superfamily (Figs 22-24)

The MHC antigens and related molecules are members of the IgSF family, their membrane proximal domains being in the IgSF C1 set. HLA class II molecules, including the $\alpha 1$ and $\alpha 2$ domains of MHC class II and the $\alpha 1$ and $\beta 1$ domains of the Class II α and β chains, show similarity to IgSF sequences⁶⁸ and the Class I domains are independent structural units as shown in Fig. 22³⁴. The Class I molecules show weak sequence similarity to each other and form a similar fold containing a

IgSF by virtue of
however, their NH₂-
class I heavy chain
how no sequence
known to form an
1 and $\alpha 2$ domains
fold containing a

together might be ca-
related to the classic
organization, includi
with a class II-like o
in sequence to MH



HLA A2 α2

Figure 24. [opposite] MHC I $\alpha 2$ set superfamily domains. Residues identical in three or more of the sequences are boxed. The positions of the beta strands (β), alpha helices (α) determined for the structure of the human HLA Class I are shown above the sequences. The sequences of the following proteins are from the Swissprot database and the database accession number and residue numbers are given in brackets. MHC Class I, human HLA Class I A-2 α precursor (P01892, 115-203); CD1A, human CD1A antigen precursor (P06126, 109-199); Fc γ rat, rat gut Fc receptor precursor (P13599, 110-199); HCMV, human cytomegalovirus glycoprotein H301 precursor (P08560, 112-210); Class II A β , mouse H2 Class I A β chain precursor (P14483, 32-122); Class DQ (3) β , human MHC Class II DQ (3) β chain precursor (P06126, 109-199).

23 MHC class I α 1 set SF domains

23. MHC class Iα1 set SF domains

	β												β												β																				
MHC Class I	G	S	H	S	R	V	F	T	S	-	V	S	R	P	G	R	G	E	P	R	F	I	A	V	G	V	D	-	D	T	Q	F	V	R	F	D	S	D	A	A	S	Q	R		
Cα1A	S	F	H	V	T	M	I	A	S	F	-	N	H	S	W	K	Q	N	L	V	S	G	W	L	S	L	Q	T	H	T	W	D	S	N	S	T	I	-	-	-	-	-	-	-	
CD8	P	R	L	P	L	M	Y	H	L	A	-	V	S	D	L	S	T	G	L	P	[S]	F	W	A	T	G	W	L	-	G	A	Q	Y	L	T	Y	N	N	L	-	-	-	-	-	-
FCR	G	M	H	V	L	R	V	G	T	G	I	F	D	-	D	T	S	H	M	T	L	T	V	V	G	I	F	D	G	Q	H	F	F	T	Y	H	V	-	-	-	-	-	-	-	
HCMV	G	M	H	V	L	R	V	G	T	G	I	F	D	-	D	T	S	H	M	T	L	T	V	V	G	I	F	D	G	Q	H	F	F	T	Y	H	V	-	-	-	-	-	-	-	

together might be called the MHC superfamily. There are numerous sequences related to the classical MHC antigens and these show a Class I type structural organization, including the binding of β 2-microglobulin, with no examples so far with a class II-like organization. The Qa and Tla antigens of mice are very similar in sequence to MHC Class I antigens whereas the human CD1 antigens show

sequences related to the classical MHC antigens and these show a Class I type structural organization, including the binding of β 2-microglobulin, with no examples so far with a class II-like organization. The Qa and Tla antigens of mice are very similar in sequence to MHC Class I antigens whereas the human CD1 antigens show

Figure 23. MHC class Ia1 set SF domains

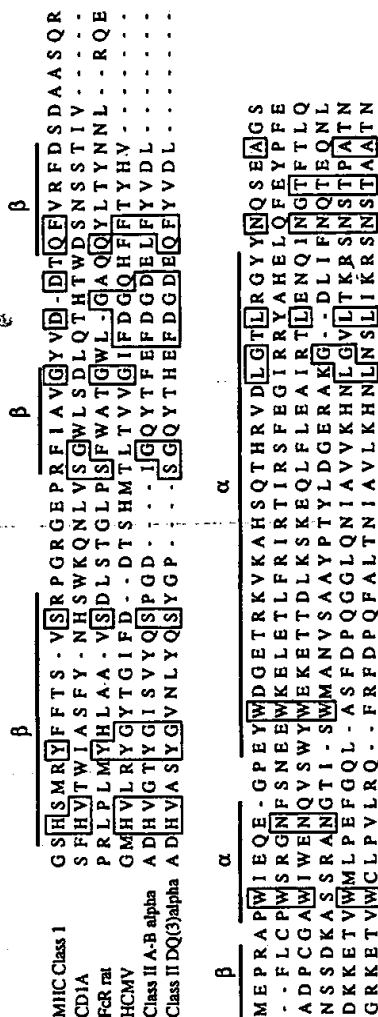
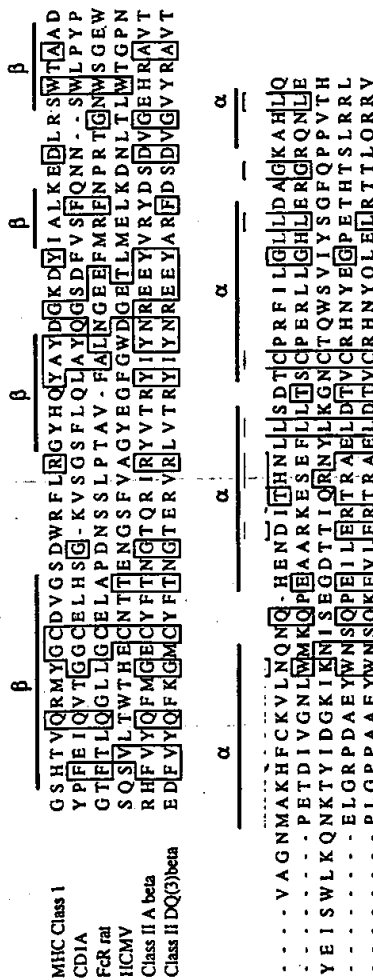


Figure 24. MHC class Ia2 set SF domains



sequence identity only at the level of about 30%. This level of identity is also seen for an Fc receptor of rodent neonatal gut ⁷⁰ and a Class I-related molecule expressed by cytomegalovirus ⁷¹. Secondary structure predictions have been used to suggest that the 70 kD heat shock proteins (hsp70) may also have a peptide binding groove like the MHC class I antigen, however the hsp70 family of proteins show little sequence similarity to the other members shown in Figs 23 and 24 ^{111, 112}.

A more detailed discussion of MHC-related sequences can be found in refs. 34, 69, 72–74.

Nerve growth factor receptor (NGFR) superfamily (Fig. 25)

Four cysteine-rich repeats were recognized in the extracellular part of the low affinity nerve growth factor receptor (NGFR) and subsequently related sequences have been identified in a number of leucocyte cell surface antigens including CD40, MRC OX-40, CD27 and the tumour necrosis factor receptor. Figure 25 shows an alignment of some of the repeats. This repeat is unusual in that most of the NGFRSF molecules contain 3 or 4 repeats. No single NGFRSF repeat sequence has been found and the repeat has not been associated with any other domain types. The gene structure of the NGFR shows that the repeat is not coded for by a single exon ⁷⁵ so it is unlikely that this repeat arose by gene duplication of exons encoding single repeats. A primordial gene with four repeats may have evolved by unequal crossing-over during recombination and this gene probably gave rise to all known members of the NGFR superfamily by duplication and divergence.

The rhodopsin superfamily (Fig. 26)

The members of this large superfamily of more than 50 proteins are characterized by the presence of seven hydrophobic membrane-spanning sequences and are reviewed in refs 76, 77. The proteins are oriented with the NH₂-terminus on the extracellular side and the COOH-terminus on the cytoplasmic side of the plasma membrane. Several names have been used to describe this superfamily, such as G protein coupled receptor superfamily, 7TMS (7-transmembrane) and rhodopsin superfamily. We have chosen the term rhodopsin superfamily as this was the first and best characterized member of this superfamily and does not imply any functional association which might later be shown to be inappropriate. The sequence conservation is highest in the potential transmembrane segments, with most diversity in the NH₂- and COOH-termini and the cytoplasmic loop between segments 5 and 6. Most members of the rhodopsinSF have been shown to couple to various G-proteins. Experiments using chimeric proteins have shown that the sequences contributing to G-protein attachment are found in transmembrane segments 5 and 6 and the cytoplasmic loop between them. A subset of closely related rhodopsinSF members is found on leucocytes and includes the C5aR fMLPR and the IL8R (see the entries in Section II and Fig. 26).

The scavenger receptor (scavengerR) superfamily (Fig. 27)

Three domains with sequence similarities were identified in the extracellular region of the CD5 antigen and later these domains were detected in macrophage scavenger receptors, the complement control protein factor I, the CD6 antigen and the speract receptor protein present in sea urchins ⁷⁸. We call this the scavengerR superfamily since these molecules are the first of this superfamily with which

Figure 25. Nerve growth factor receptor (NGFR) superfamily repeats. Residues identical in five or more sequences are boxed. The asterisks mark the positions of the conserved residues that are marked on the domain organization figures in the entries in Section II. The sequences of the following proteins are from the Swissprot database unless otherwise indicated and the database accession number and residue numbers are given in brackets. The sequences are contiguous over the four repeats except for OX-40 which contains a short sequence in place of the third repeat. OX-40, rat MRC OX-40 antigen precursor (P15725, 25–102, 124–164); TNFR1, human tumour necrosis factor receptor precursor I (P19438, 43–196); TNFR2, human tumour necrosis factor receptor precursor II (P20333, 39–201); NGFR, rat nerve growth factor receptor precursor (P07174, 32–190); CD40, human

Figure 25. Nerve growth factor receptor (NGFR) superfamily repeats. Residues identical in five or more sequences are boxed. The asterisks mark the positions of the conserved residues that are marked on the domain organization figures in the entries in Section II. The sequences of the following proteins are from the Swissprot database unless otherwise indicated and the database accession number and residue numbers are given in brackets. The sequences are contiguous over the four repeats except for OX-40 which contains a short sequence in place of the third repeat. OX-40, rat MRC OX-40 antigen precursor (P15725, 25-102, 124-164); TNFRI, human tumour necrosis factor receptor precursor I (P19438, 43-196); TNFRII, human tumour necrosis factor receptor precursor II (P20333, 39-201); NGFR, rat nerve growth factor receptor precursor (P07174, 32-190); CD40, human CDw40 antigen precursor (PIR:S04460, 25-187).

OX40 (1)	N	C	V	K	D	T	T	P	S	...	G	H	K	C	...	C	...	R	E	C	Q	P	...	G	H	G	M	V	S	R	C	D	H	...	T	R	D	T	V	C	H								
TNFRI(1)	V	C	P	Q	G	K	Y	I	H	P	Q	N	...	N	S	I	C	...	T	K	C	H	K	...	G	T	Y	L	Y	N	D	C	P	G	Q	D	T	D	C	R									
TNFRII(1)	T	C	R	L	R	E	Y	Y	D	Q	T	...	A	Q	M	C	...	S	K	C	S	P	...	G	H	A	K	V	F	C	T	K	...	T	S	D	T	V	C	V									
NGFR (1)	T	C	S	T	G	L	Y	T	H	...	S	G	E	C	...	C	...	K	A	C	N	L	...	G	E	G	V	A	Q	P	C	G	...	A	N	Q	T	V	C	E									
CD40 (1)	A	C	R	E	K	Q	Y	L	I	...	N	S	Q	C	...	C	...	S	L	C	Q	P	...	G	Q	K	L	V	S	D	C	T	E	...	F	T	E	T	E	C	L								
OX40 (2)	P	C	E	P	G	F	Y	N	E	A	V	N	Y	...	D	T	C	K	Q	C	...	T	Q	C	N	H	R	S	G	S	E	L	K	Q	N	C	T	P	...	T	E	D	T	V	C				
TNFRI(2)	E	C	E	S	G	S	F	T	A	S	E	N	H	...	L	R	H	C	L	S	C	...	S	K	C	R	K	E	M	G	Q	V	E	I	S	S	C	T	V	...	D	R	D	T	V	C			
TNFRII(2)	S	C	E	D	S	T	Y	T	Q	L	W	N	W	...	V	P	E	C	L	S	C	...	G	S	R	C	S	...	D	Q	V	E	T	Q	A	C	T	R	...	E	Q	N	R	I	C				
NGFR (2)	P	C	L	D	N	V	T	F	S	D	V	S	A	T	E	P	C	K	P	C	...	T	E	C	L	G	...	L	Q	S	M	S	A	P	C	V	E	...	A	D	D	A	V	C					
CD40 (2)	P	C	G	E	S	E	F	L	D	T	W	N	R	E	...	T	H	C	H	Q	H	...	K	Y	C	D	P	N	L	G	L	R	V	Q	K	G	I	S	...	E	T	D	T	I	C				
TNFRI(3)	G	C	R	K	N	Q	Y	R	H	Y	W	S	E	N	L	F	Q	C	F	N	C	...	S	L	C	L	N	...	G	T	V	H	L	S	...	C	Q	E	...	K	Q	N	T	V	C	T			
TNFRII(3)	T	C	R	P	G	W	Y	C	A	L	S	K	Q	...	E	G	C	R	L	C	...	A	P	L	R	K	C	R	P	...	G	F	G	V	A	R	P	G	T	E	...	T	S	D	V	V	C	K	
NGFR (3)	R	C	A	Y	G	Y	Y	Q	D	E	E	T	...	G	H	C	E	A	C	...	S	V	C	E	V	...	G	S	C	L	V	F	S	C	D	...	K	Q	N	T	V	C	E						
CD40 (3)	T	C	E	E	G	W	H	C	T	S	E	...	A	C	E	S	C	V	L	H	R	S	C	S	P	...	G	F	G	V	K	Q	I	A	T	G	...	V	S	D	T	T	C	E					
OX40 (4)	P	C	P	P	G	H	F	S	P	G	S	N	Q	...	A	C	K	P	W	...	T	N	C	T	L	...	S	G	K	Q	I	R	H	P	A	S	N	...	S	L	D	T	V	C	E				
TNFRI(4)	C	H	A	G	P	F	L	R	E	N	...	E	C	V	S	C	...	S	N	C	K	K	...	S	N	C	K	K	...	S	L	E	C	T	K	...	S	M	D	A	V	C	T						
TNFRII(4)	P	C	A	P	G	T	F	S	N	T	T	S	...	T	D	I	C	R	P	H	...	Q	I	C	N	V	...	V	A	I	P	G	N	A	...	S	M	D	A	V	C	T							
CD40 (4)	P	C	P	V	G	F	E	S	N	V	S	S	A	...	F	E	K	C	H	P	W	...	T	S	C	I	E	T	...	K	D	L	V	V	Q	I	A	G	T	I	N	...	K	T	I	V	V	V	L

Figure 26. Rhodopsin superfamily

IL8R	MSNITDPQMWDFDDLN....FTOMPADBDYSPCMLBETLTNKYVVIAYAVFVLSLVLGNLSVILYSVGR
C5aR	...MNSFNYYTTPDGYHDDTDLNTPVDKT...SNTLRVPDILALVIFAVFVLGVGLGNALVWVWAF-BAXR
MLPR	...MEINSLPTNISGGTPAYSGYLFLLIITLVAVFVGLGNGLVWVAGF-BMTH
Rhodopsin	...MNGTEGPNFYVFNATGVVRSFPEYQYLAEPQFMSLAAYMFLLI-VLGFVNLTLVTVQHKLR
NeurokininR	-MGTCDIVTEANISSGPESNTGIT.....AFSMPSWQLALWAPAYLVVAVVGNALVWVILAHRRMR
DopaRMRTLNTSAMDGTOLVVERDFSVRLTACFLSLIISTLEGNLTVCAAVIRFHLR
I	
II	
III	
IV	
V	
VI	
VII	

Figure 26. (above) Rhodopsin superfamily. Residues identical in four or more sequences are boxed. The bars over the sequences indicate the transmembrane regions. The sequences of the following proteins are from the Swissprot database unless otherwise indicated and the database accession number. IL8R, human high affinity IL8 receptor 109; C5aR, human C5a anaphylatoxin chemotactic receptor (P21730); fMLPR, human fMet-Leu-Phe receptor (P21462); Rhodopsin, human rhodopsin (P08100); NeurokininR, human neurokinin A receptor (P21452); DopaR, human D(1) dopamine receptor (P21728).

Figure 27. (below) Scavenger receptor superfamily domain. Residues identical in four or more sequences are boxed. The asterisks mark the positions of the conserved residues that are marked on the domain organization figures in the entries in Section II. The sequences of the following proteins are from the Swissprot database unless otherwise indicated and the database accession

Figure 26. (above) Rhodopsin superfamily. Residues identical in four or more sequences are boxed. The bars over the sequences indicate the transmembrane regions. The sequences of the following proteins are from the Swissprot database unless otherwise indicated and the database accession number. IL8R, human high affinity IL8 receptor 109; C5aR, human C5a anaphylatoxin chemotactic receptor (P21730); fMLPR, human fMet-Leu-Phe receptor (P21462); Rhodopsin, human rhodopsin (P08100); NeurokininR, human neurokinin A receptor (P21452); DopaR, human D(1) dopamine receptor (P21728).

Figure 27. (below) Scavenger receptor superfamily domain. Residues identical in four or more sequences are boxed. The asterisks mark the positions of the conserved residues that are marked on the domain organization figures in the entries in Section II. The sequences of the following proteins are from the Swissprot database unless otherwise indicated and the database accession number and residue numbers are given in brackets. SREC, sea urchin egg peptide speract precursor (P16264, d1 38–145, d2 148–258, d3 259–367, d4 377–486); CD5, human CD5 antigen (P06127, d1 30–134, d2 156–269, d3 271–369); CFAL, human complement factor I precursor (P05156, 109–216); SC5V, human scavenger receptor type I (P21757, 345–451). Alignments are from 20 residues from the conserved glycine to the COOH-terminus of the scavenger receptor.

[illegible]

1
 2
 3
 4
 5
 6
 7
 8
 9
 10
 11
 12
 13
 14
 15
 16
 17
 18
 19
 20
 21
 22
 23
 24
 25
 26
 27
 28
 29
 30
 31
 32
 33
 34
 35
 36
 37
 38
 39
 40
 41
 42
 43
 44
 45
 46
 47
 48
 49
 50
 51
 52
 53
 54
 55
 56
 57
 58
 59
 60
 61
 62
 63
 64
 65
 66
 67
 68
 69
 70
 71
 72
 73
 74
 75
 76
 77
 78
 79
 80
 81
 82
 83
 84
 85
 86
 87
 88
 89
 90
 91
 92
 93
 94
 95
 96
 97
 98
 99
 100
 101
 102
 103
 104
 105
 106
 107
 108
 109
 110
 111
 112
 113
 114
 115
 116
 117
 118
 119
 120
 121
 122
 123
 124
 125
 126
 127
 128
 129
 130
 131
 132
 133
 134
 135
 136
 137
 138
 139
 140
 141
 142
 143
 144
 145
 146
 147
 148
 149
 150
 151
 152
 153
 154
 155
 156
 157
 158
 159
 160
 161
 162
 163
 164
 165
 166
 167
 168
 169
 170
 171
 172
 173
 174
 175
 176
 177
 178
 179
 180
 181
 182
 183
 184
 185
 186
 187
 188
 189
 190
 191
 192
 193
 194
 195
 196
 197
 198
 199
 200
 201
 202
 203
 204
 205
 206
 207
 208
 209
 210
 211
 212
 213
 214
 215
 216
 217
 218
 219
 220
 221
 222
 223
 224
 225
 226
 227
 228
 229
 230
 231
 232
 233
 234
 235
 236
 237
 238
 239
 240
 241
 242
 243
 244
 245
 246
 247
 248
 249
 250
 251
 252
 253
 254
 255
 256
 257
 258
 259
 260
 261
 262
 263
 264
 265
 266
 267
 268
 269
 270
 271
 272
 273
 274
 275
 276
 277
 278
 279
 280
 281
 282
 283
 284
 285
 286
 287
 288
 289
 290
 291
 292
 293
 294
 295
 296
 297
 298
 299
 300
 301
 302
 303
 304
 305
 306
 307
 308
 309
 310
 311
 312
 313
 314
 315
 316
 317
 318
 319
 320
 321
 322
 323
 324
 325
 326
 327
 328
 329
 330
 331
 332
 333
 334
 335
 336
 337
 338
 339
 340
 341
 342
 343
 344
 345
 346
 347
 348
 349
 350
 351
 352
 353
 354
 355
 356
 357
 358
 359
 360
 361
 362
 363
 364
 365
 366
 367
 368
 369
 370
 371
 372
 373
 374
 375
 376
 377
 378
 379
 380
 381
 382
 383
 384
 385
 386
 387
 388
 389
 390
 391
 392
 393
 394
 395
 396
 397
 398
 399
 400
 401
 402
 403
 404
 405
 406
 407
 408
 409
 410
 411
 412
 413
 414
 415
 416
 417
 418
 419
 420
 421
 422
 423
 424
 425
 426
 427
 428
 429
 430
 431
 432
 433
 434
 435
 436
 437
 438
 439
 440
 441
 442
 443
 444
 445
 446
 447
 448
 449
 450
 451
 452
 453
 454
 455
 456
 457
 458
 459
 460
 461
 462
 463
 464
 465
 466
 467
 468
 469
 470
 471
 472
 473
 474
 475
 476
 477
 478
 479
 480
 481
 482
 483
 484
 485
 486
 487
 488
 489
 490
 491
 492
 493
 494
 495
 496
 497
 498
 499
 500
 501
 502
 503
 504
 505
 506
 507
 508
 509
 510
 511
 512
 513
 514
 515
 516
 517
 518
 519
 520
 521
 522
 523
 524
 525

clear functional activity has been associated. However, some ligands will bind to both scavenger receptors I and II but the latter lacks this type of domain so the functional involvement of this domain remains to be resolved ⁷⁹. Initially it was argued that the CD5 antigen domains were related to IgSF domains ⁸⁰. However, this contention was not supported by ALIGN analysis as described above. Subsequently it was suggested that the CD5 domains were related in sequence to the domains of the PapD bacterial protein ¹¹ but again this was not supported by a detailed analysis including CD5 domains plus numerous other scavengerRSF domains. It is now clear that there is a separate superfamily of proteins containing scavengerRSF domains and alignments for this superfamily are shown in Fig. 27. No tertiary structure data are available yet for these domains.

Signal transduction sequence motifs (Fig. 28)

The signal transduction sequence motif shown in alignments in Fig. 28 is present in the cytoplasmic regions of several membrane proteins present in the antigen receptor complexes on B cells, T cells, and the IgE receptor on mast cells ⁸¹. This motif is also found in the CD5 antigen cytoplasmic domain (Beyers, Spruyt and Williams, unpublished). The CD3 ζ chain is unusual in that it has three motifs whereas the others have only one. One common feature of these molecules is that they are components of membrane complexes which, when cross-linked, give signals that lead to cell activation. This results in cell proliferation in the case of the antigen receptors and to degranulation of mast cells. Cross-linking of CD3 ϵ by

Human CD3 gamma	D	K	Q	T	L	L	P	N	D	Q	L	Y	Q	P	L	K	D	R	E	D	D	Q	-	Y	S	H	L	Q	G	N	Q	L	R	R	N
Human CD3 delta	D	T	Q	A	L	L	R	N	D	Q	V	Y	Q	P	L	R	D	R	D	D	A	Q	-	Y	S	H	L	G	G	N	W	A	R	N	K
Mouse CD3 epsilon	K	E	R	P	P	P	V	P	N	P	D	Y	E	P	I	R	K	G	Q	R	D	L	-	Y	S	G	L	-	-	-	-	-	-	-	
Human CD3 zeta (1)	P	P	A	Y	Q	Q	Q	N	Q	L	Y	N	E	L	N	L	G	R	E	E	-	Y	D	V	L	D	K	R	R	G	R	D	P		
Human CD3 zeta (2)	K	P	R	R	K	N	P	Q	E	G	L	Y	N	E	L	Q	K	D	K	M	A	E	A	Y	S	E	I	G	M	K	G	-	-	-	
Human CD3 zeta (3)	E	R	R	R	G	K	G	H	D	G	L	Y	Q	G	L	S	T	A	T	K	D	T	-	Y	D	A	L	H	M	Q	A	L	P	P	
Mouse MB1	D	M	P	D	D	Y	E	D	E	N	L	Y	E	G	L	N	L	D	D	C	S	M	-	Y	E	D	I	S	R	G	L	Q	G	T	Y
Mouse B29	D	G	K	A	G	M	E	D	H	T	Y	E	G	L	N	I	D	Q	T	A	T	-	Y	E	D	I	V	T	L	R	T	G	E	V	
Rat Fc epsilon R beta chain	F	E	R	S	K	V	P	D	D	R	L	Y	E	E	L	H	V	S	P	I	-	-	Y	S	A	L	E	D	T	R	E	A	S	A	
Rat Fc epsilon R gamma chain	D	I	A	S	R	E	K	S	D	A	V	Y	T	G	L	N	T	R	N	Q	E	T	-	Y	E	T	L	K	H	E	K	P	P	Q	
Human CD5	E	N	P	T	A	S	H	V	D	N	E	Y	S	Q	P	P	R	N	S	R	L	S	A	Y	P	A	L	E	G	V	L	H	R	S	-

Figure 28. Signal transduction motifs. Residues identical in five or more sequences are boxed. The sequences of the following proteins are from the Swissprot database and the database accession number and residue numbers are given in brackets. CD3 gamma, human CD γ chain precursor (P09693, 149-181); CD3 delta, human CD3 δ chain precursor (P04234, 138-171); CD3 epsilon, mouse CD3 ϵ chain precursor (P22646, 159-184); CD3 zeta, human CD3 ζ chain precursor (P20963, 61-94, 99-130, 131-163); MB1, mouse MB-1 protein precursor (P11911, 171-204); B29, mouse B cell glycoprotein B29 precursor (P15530, 184-217); Fc epsilon R beta, rat Ig ϵ receptor β subunit (P13386, 207-239); Fc epsilon R gamma, rat Ig ϵ receptor γ subunit precursor (P20411, 54-86); CD5, human CD5 antigen precursor (P06127, 442-475).

PC1 d1	K	S	C	-	K	G	R	C
PP11	T	S	C	-	Q	G	R	C
Vitronectin	E	S	C	-	K	G	R	C
PC1 d2	W	T	C	N	K	F	R	C

Figure 29. Somatomed more sequences are bo residues that are mark Section II. The sequen database and the data brackets. PC1, mouse PP11, human placenta vitronectin precursor (

immobilized mAbs le shown that the cross chain or the CD3 ϵ , give TcR, implying that transduction mechan

Somatomedin B super
Somatomedin B is a spreading factor) by p contain two somato pyrophosphatase/alka a different region of repeat. The domain h is present in placental

Transmembrane 4 pa Chain and CD20 (Figs
The "TM4 superfamily clear sequence simila with both the NH₂- a This superfamily incl CD63 and TAPA-1. A genomic sequence of eight exons which transmembrane sequ largely compatible w molecules had a com sequence between T between TM sequer sequence length. Thi includes the N-link labelled at the cell addition, surface lab

ven in
 D3
 use CD3
 rsor
 1911,
 : Fc
 gamma,
 tigen

Protein Superfamilies and Cell Surface Molecules

PC1 d1 K S C - K G R C F E - - R T F S N C R C D A A C V S L G N C C K L D Q E T C V E P T H
 PP11 T S C - Q G R C Y E A F D K D K H H C C H C N A R C Q E F G N C C K L D E S L C S D H E V
 Vitronectin E S C - K G R C Y E G F N V D K K C C D E L C S Y Y Q C S C T D T A E C K P Q V T
 PC1 d2 W T C N K F R C G E K R L S R F V C S C A D D C K T H N D C C I N S S V C Q D K K S

Figure 29. Somatomedin B superfamily domains. Residues identical to three or more sequences are boxed. The asterisks mark the positions of the conserved residues that are marked on the domain organization figures in the series in Section II. The sequences of the following proteins are from the Swiss-Prot database and the database accession number and residue numbers are given in brackets. PC1, mouse plasma cell antigen PC1 (P06802; d1, 54-93; d2, 95-137); PP11, human placental protein precursor (P21128, 47-88); Vitronectin, human vitronectin precursor (P04004; 22-63).

immobilized mAbs leads to T cell proliferation. The use of chimeric receptors has also been shown that the cross-linking of the cytoplasmic domains of the TCR β chain or the CD3 ϵ , gives a TcR-like signal in cells lacking surface expression of the TcR, implying that this motif is involved in coupling the TcR to downstream transduction mechanisms^{82,83}.

Somatomedin B superfamily (Fig. 29)

Somatomedin B is a serum peptide derived from vitronectin (also spreading factor) by proteolysis. The plasma cell surface antigen PC-1 contains two somatomedin BSF repeats⁴¹. This glycoprotein has acid phosphatase/alkaline phosphodiesterase activity⁸⁴ but this is a different region of the molecule than that containing the somatomedin repeat. The domain has not been found on other cell surface molecules but is present in placental protein 11 (PP11)⁸⁵.

proteins has
all receptor ζ
ession of the
intracellular

called serum
was noted to
a nucleotide
associated with
tomodisin BSF
as although it-

Transmembrane 4 pass (TM4) superfamily and the relationship between FcεRIβ Chain and CD20 (Figs 30 and 31)

The "TM4 superfamily" is a term that we suggest for a new group of clear sequence similarities that are thought to traverse the lipid bilayer with both the NH₂- and COOH-termini on the cytoplasmic face of the membrane. This superfamily includes several leucocyte antigens such as CD9, CD63 and TAPA-1. Alignments for the TM4 superfamily are shown in Figure 1. The genomic sequence of the TAPA-1 antigen shows that the sequence of eight exons which do not indicate any simple correlation with transmembrane sequences⁸⁶. The intron/exon boundaries of most of the TM4 superfamily molecules are largely compatible with those of TAPA-1⁸⁷ in support of the argument that these molecules had a common ancestor in evolution. The majority of the sequence between TM4 superfamily molecules reside in the extracellular domain between TM sequences 3 and 4 where there are considerable differences in sequence length. This loop of sequence is known to be extracellular and includes the N-linked glycosylation sites and the MRC OX-44 epitope. The MRC OX-44 epitope is labelled at the cell surface maps to an Ile/Thr interchange in the extracellular domain. In addition, surface labelling studies on TAPA-1 support an extracellular

proteins with
er four times
e membrane.
CD37, CD53,
a Fig. 30. The
is coded by
the proposed
se CD53 are
ent that these
differences in
acellular loop
differences in
ar because it
e that can be
region⁸⁷. In
ur localization

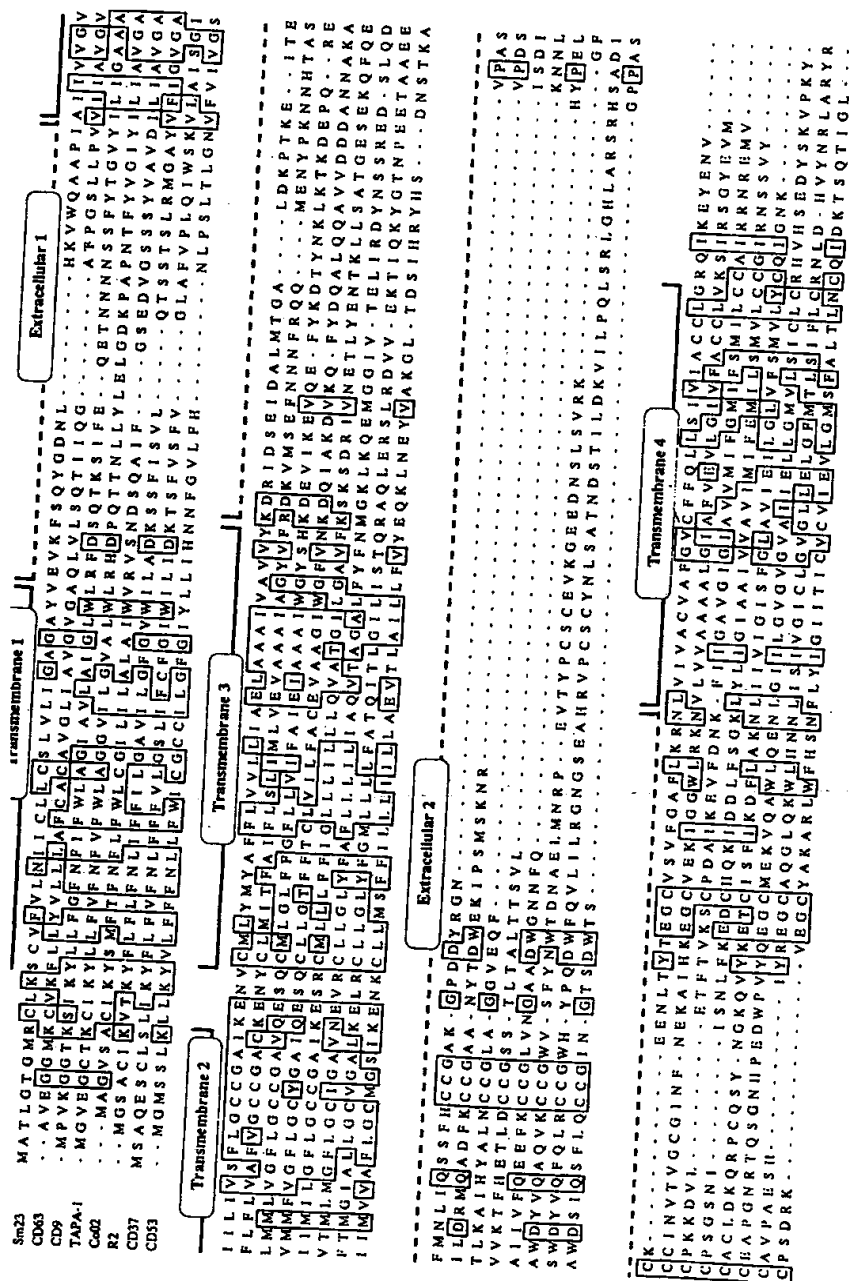
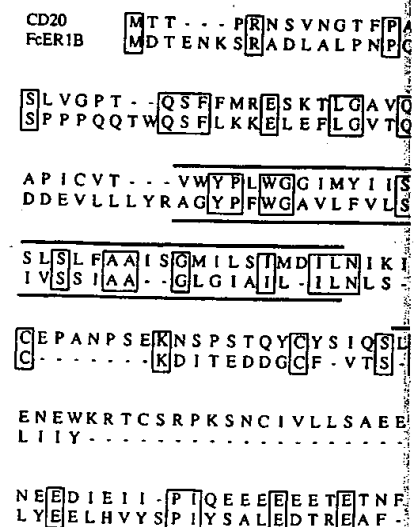


Figure 30. (opposite) Transmembrane 4 in four or more sequences are boxed. T from the Swissprot database or the reference number is given in brackets (the complete Schistosoma mansoni protein Sm23 (P08962); CD9, human antigen ME491 (P08962); CD9, human TAPA-1 antigen (P18582); Co029, human (P19075); R2, human R2 antigen 110, Co human CD53 antigen (P19397).

Figure 31. (below) Alignment of CD20 between the two sequences are boxed. transmembrane sequences are indicated. The similarities are mostly within or at fall off towards the COOH-terminus. T and an ALIGN score of 6.1 SD was obtained from the Swissprot database accession numbers: chain (P20490).



for this loop ⁸⁸. There is no known function although antibodies against members proliferate of leucocytes ⁸⁹.

There are other leucocyte proteins that membrane four times, including CD20. These show sequence similarity to the TM4 superfamily. The chain show clear sequence similarities in transmembrane regions of these molecules.

Figure 31. (below) Alignment of CD20 and FcεR1 β chain. Residues identical between the two sequences are boxed. The possible positions of the four transmembrane sequences are indicated with a bar above or below the sequences. The similarities are mostly within or around the transmembrane sequences and fall off towards the COOH-terminus. There are many conservative substitutions and an ALIGN score of 6.1 SD was obtained for the full sequences shown. The Swissprot database accession numbers are; human CD20 (P22836), mouse FcεR1 β chain (P20490).

CD20
FcER1B

MTT---PRNSVNGTFPAEPMKGP-IAM-QSGP--KPLFRMS
MDTENKSRADLALPNPQESPSAPDIELLEASPPAKALPEKPA

SLVQPT--QSFFMR[ESKTLGAVQIMNGLFHIALGGLL--MIPAGIY
SPPPPQQTWQSFLKKELEFLGVTVQLVGLICLCFGTVVVCSTLQTSDF

APICVT--VWYPLWGGIMYIISGSLLAATEKNSRKCLVKGKMMIN
DDEVLLLYRAGYPFWGAVLFVLSGFLSIMS[ERKNTLYLVRGSLGAN

SLSLFAAISGMILSTMDILNLIKISHFLKMESLNFIRAHTPYININYN
IVSSIAA--GLGLAIIILNLS-----NNSAYMNY-

CEPANPSEKNSPSTQYCYSIQSLFLGLILSVMLIFAFFQELVIAGIV
C-----KDITEDDGCF-VTS-FITELVLVLMFLFLTLAFCSAFVL

ENEWKRTCSPKSN[NCIVLLSABEKKETIEIKEEVVGLTETETSSQPK
LII-----RIGQEFB-RSKV-----PDDR

NEEDIEII-PIQEEEBEETETNFPEPPQDCQESSPIENDSSP
LYEELHVYSPILYSAL[EDTREAFAF--SAPVVS----

for this loop ⁸⁸. There is no known function for any of this family of proteins although antibodies against members of this family do have effects on the proliferation of leucocytes ⁸⁹.

There are other leucocyte proteins that are predicted to traverse the plasma membrane four times, including CD20 and the FcεRIβ chain, but these do not show sequence similarity to the TM4 superfamily. However, CD20 and the FcεRIβ chain show clear sequence similarities to each other⁹⁰ in three of the four transmembrane regions of these molecules as shown in Fig. 31, and their genes are

very closely linked on mouse chromosome 19⁹⁰. Thus these two sequences should be considered as founder members of a new superfamily and perhaps for now this could be referred to as the FcεRIβ superfamily. There are data to show that CD20 is a Ca²⁺ channel⁹¹ and it would be interesting to know if this is also the case for FcεRIβ chain and members of the TM4SF too.

The tyrosine kinase superfamily (Fig. 32)

Tyrosine kinase domains are found in the cytoplasm and two groups can be distinguished: receptor tyrosine kinases which are transmembrane proteins, and nonreceptor tyrosine kinases which are located in the cytoplasm. The non-receptor group of kinases includes members of the src family, all of which are anchored to the inner leaflet of the plasma membrane with a myristate moiety. On activation they phosphorylate Tyr residues on their own cytoplasmic domains or on other proteins in the cytoplasm and this is believed to be one of the key early events in signal transduction pathways after ligand recognition. In leucocytes the best studied example is p56^{lck} which associates with the cytoplasmic domains of CD4 and CD8 and regulates signal transduction by these molecules ⁹². Other examples are *fyn* which associates with the T cell receptor complex ⁹³, and *lyn*, *fyn*, and *blk*, which couple to the membrane Ig complex of B cells ^{94, 95}.

Receptor tyrosine kinases are expressed on a wide variety of cells and examples include the PDGF receptor and EGF receptor. When these receptors bind their natural ligands they oligomerize and the cytoplasmic tyrosine kinase domains become activated and autophosphorylated. This leads to the phosphorylation and activation of various intracellular substrates including phospholipase C γ , phosphatidylinositol 3-kinase and the c-raf serine kinase. These effector molecules concomitantly associate with the activated receptor kinases^{96,97}.

Tyrosine kinase domains consist of about 260–360 amino acids. The difference in size is due to insertion of a “kinase insert domain” of about 70–100 amino acids in certain receptor kinases, including the platelet derived growth factor receptor, M-CSFR, and c-kit kinases. These insert regions appear to regulate the interaction of the kinase with certain cellular substrates/effector molecules^{96,97}. The tyrosine kinase domain of a particular molecule is particularly well-conserved across species and the identities between molecules within the superfamily are about 40%, as illustrated in Fig. 32. This is much higher than for many of the superfamilies with domains that are found at the cell surface.

Kinase domains are not conserved uniformly, but consist of 11 highly conserved subdomains (I–XI) separated by regions of lower conservation⁹⁸. Subdomain I contains the Gly-X-Gly-X-X-Gly consensus which forms part of the binding site for ATP. Subdomain II contains an invariant lysine, which appears to be directly involved in the phosphotransfer reaction. Subdomain VIII contains a Pro-Ile-/Val-/Arg-Trp-Thr/Met-Ala-Pro-Glu consensus which is characteristic of the tyrosine kinases. In the serine/threonine kinases the consensus is Gly-Thr/Ser-X-X-Tyr/Phe-Ala-Pro-Glu.

Tyrosine kinases have been the subject of extensive study in nonlymphoid cells and the subject has been reviewed in depth 96, 98, 99.

e phospho-tyrosine phosphatase (PTPase) superfamily (Fig. 33)

e PTPase superfamily of integral membrane proteins was discovered when two cytoplasmic homology units of the CD45 antigen¹⁰⁰ were matched with the

	I	II	III	IV
SRC	ES	LR	LV	KL
LCK	ET	KL	VR	KL
MCSEIR	NN	QF	GT	KL
KIT	NR	LS	FG	KL
EGPR	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR	LV	KL
	ET	KL	VR	KL
	NN	QF	GT	KL
	NR	LS	FG	KL
	TE	FK	IK	KL
	ES	LR		

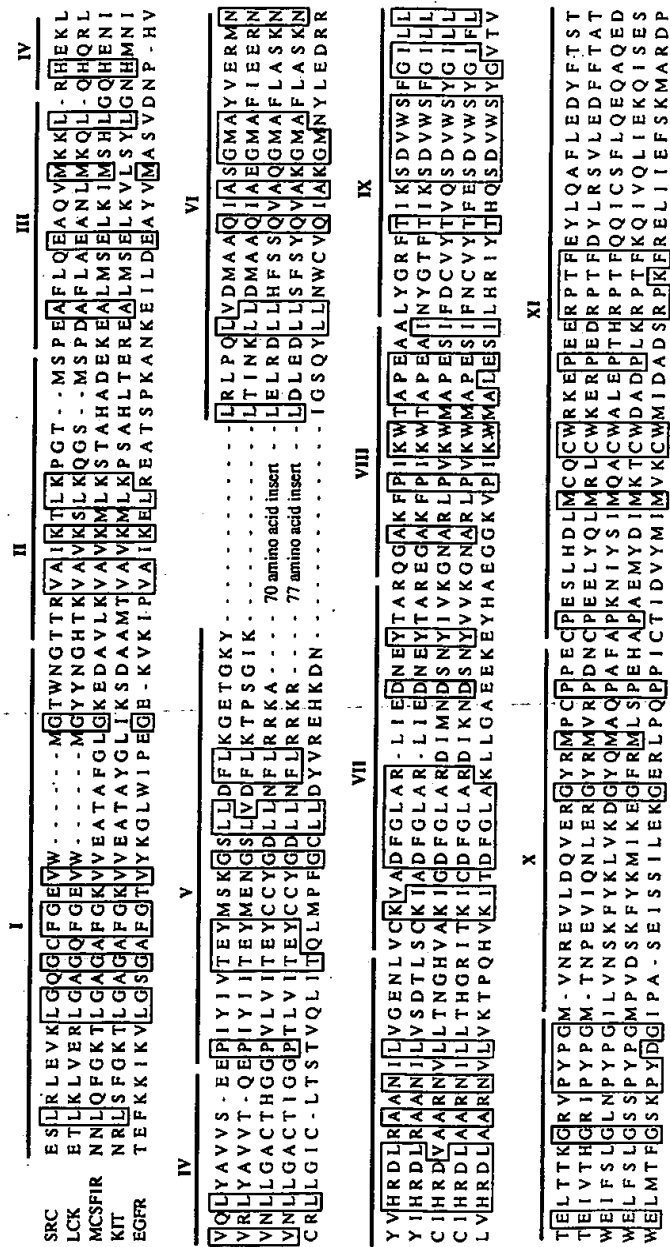


Figure 32. Tyrosine kinase superfamily. Residues identical in four or more sequences are boxed. The bars above the sequences represent the subdomains defined in this superfamily and the numbering is as in ref. 98. The sequences of the following proteins are from the Swissprot database and the database accession number and residue numbers are given in brackets. SRC, human src proto-oncogene tyrosine kinase (P12931, 268-526); LCK, human T cell specific tyrosine kinase (P06239, 243-501); MCSFR, human macrophage colony stimulating factor receptor precursor (P07333, 580-917); KIT, human kit proto-oncogene precursor (P10721, 587-931); EGFR, human EGF receptor precursor (P00533, 710-975).

sequence of a placental cytoplasmic phospho-tyrosine phosphatase^{101,102}. Subsequently PTPase activity has been shown for the membrane proximal cytoplasmic domain of CD45 but as yet not for the COOH-terminal domain¹⁰³. Subsequently, other sequences have been identified with similarities to these sequences by cross-hybridization with cDNA probes and these include the LAR protein that is a cell surface protein with three IgSF domains and eight FN type III SF domains on the extracellular side, although this protein is not expressed widely on leucocytes¹⁰⁴. Sequences with many similarities to the PTPaseSF domains have been identified in *Drosophila*¹⁰⁵. Some examples of the sequences are shown in Fig. 33. The second domain in CD45 is unusual in comparison with all other members of the PTPaseSF in that it contains an insertion of 19 amino acids with a very high content of acidic and Ser residues. The Ser residues may be phosphorylated by Ser kinases to produce an extremely negatively charged region of sequence.

The only complete genomic structure at present in the PTPaseSF is for mouse CD45¹⁰⁶. This shows that the region illustrated in Fig. 33 is encoded by 6 or 8 exons for each domain. However, the ends of the domains as defined from the sequence similarities do not correspond to the ends of exons with the same phase of intron/exon boundaries. The genetic origin of these domains is unclear.

References

- 1 Dayhoff, M.O. et al. (1983) *Methods Enzymol.* 91, 524-45.
- 2 Williams, A.F. and Barclay, A.N. (1988) *Annu. Rev. Immunol.* 6, 381-405.
- 3 Baron, M. et al. (1991) *Trends Biochem. Sci.* 16, 13-17.
- 4 Pearson, W.R. and Lipman, D.J. (1988) *Proc. Natl Acad. Sci. USA* 85, 2444-2448.
- 5 Patthy, L. (1990) *Cell* 61, 13-14.
- 6 Williams, A.F. (1987) *Immunol. Today* 8, 298-303.
- 7 Williams, A.F. et al. (1989) *Cold. Spring. Harb. Symp. Quant. Biol.* LIV, 637-647.
- 8 Ryu, S.E. et al. (1990) *Nature (London)* 348, 419-426.
- 9 Wang, J. et al. (1990) *Nature (London)* 348, 411-418.
- 10 Driscoll, P.C. et al. (1991) *Nature (London)* 353, 762-765.
- 11 Holmgren, A. and Branden, C.I. (1989) *Nature (London)* 342, 248-251.
- 12 Amzel, L.M. and Poljak, R.J. (1979) *Annu. Rev. Biochem.* 48, 961-997.
- 13 de Vos, A.M. et al. (1992) *Science* 255, 306-312.
- 14 Baron, M. et al. (1992) *Biochemistry* 31, 2068-2073.
- 15 Patthy, L. (1987) *J. Mol. Biol.* 198, 567-577.
- 16 Sharp, P.A. (1981) *Cell* 23, 643-646.
- 17 Barton, G.J. and Sternberg, M.J. (1987) *J. Mol. Biol.* 198, 327-337.
- 18 Devereux, J. et al. (1984) *Nucl. Acids Res.* 12, 387-395.
- 19 Reid, K.B. and Day, A.J. (1989) *Immunol. Today* 10, 177-180.
- 20 Barlow, P.N. et al. (1991) *Biochemistry* 30, 997-1004.
- 21 Idzerda, R.L. et al. (1990) *J. Exp. Med.* 171, 861-873.
- 22 Cosman, D. et al. (1990) *Trends Biochem. Sci.* 15, 265-270.
- 23 Goodwin, R.G. et al. (1990) *Cell* 60, 941-951.
- 24 Bazan, J.F. (1990) *Proc. Natl Acad. Sci. USA* 87, 6934-6938.
- 25 Fukunaga, R. et al. (1991) *EMBO J.* 10, 2855-2865.
- 26 Tappin, M.J. et al. (1989) *Eur. J. Biochem* 179, 629-637.
- 27 Cooke, R.M. et al. (1987) *Nature (London)* 327, 339-341.
- 28 Handford, P.A. et al. (1990) *EMBO J.* 9, 475-480.

- 29 Rebay, I., et al. (1991) *Cell* 67, 687-699.
- 30 Constantine, K.L. et al. (1991) *Biochemistry* 30, 1663-1672.
- 31 Benian, G.M. et al. (1989) *Nature (London)* 342, 45-50.
- 32 Labeit, S. et al. (1990) *Nature (London)* 345, 273-276.
- 33 Suzuki, S. and Naitoh, Y. (1990) *EMBO J.* 9, 757-763.
- 34 Bjorkman, P.J. et al. (1987) *Nature (London)* 329, 506-512.
- 35 Leahy, D.J. et al. (1992) *Cell* 68, 1145-1162.
- 36 Piggott, R. and Power, C., (1992) *The Adhesion Molecule FactsBook*, Academic Press, London.
- 37 Erle, D.J. et al. (1991) *J. Biol. Chem.* 266, 11009-11016.
- 38 Takada, Y. and Hemler, M.E. (1989) *J. Cell Biol.* 109, 397-407.
- 39 Hemler, M.E. (1990) *Annu. Rev. Immunol.* 8, 365-400.
- 40 Drickamer, K. (1988) *J. Biol. Chem.* 263, 9557-9560.
- 41 Patthy, L. (1988) *J. Mol. Biol.* 202, 689-696.
- 42 Lasky, L.A. et al. (1989) *Cell* 56, 1045-1055.
- 43 Hoyle, G.W. and Hill, R.L. (1991) *J. Biol. Chem.* 266, 1850-1857.
- 44 Doege, K. et al. (1987) *J. Biol. Chem.* 262, 17757-17767.
- 45 Halberg, D.F. et al. (1988) *J. Biol. Chem.* 263, 9486-9490.
- 46 Collins, T. et al. (1991) *J. Biol. Chem.* 266, 2466-2473.
- 47 Johnston, G.I. et al. (1990) *J. Biol. Chem.* 265, 21381-21385.
- 48 Ezekowitz, R.A. et al. (1990) *J. Exp. Med.* 172, 1785-1794.
- 49 Weis, W.I. et al. (1992) *Science* 254, 1608-1615.
- 50 Frigeri, L.G. et al. (1990) *J. Biol. Chem.* 265, 20763-20769.
- 51 Hickey, M.J. et al. (1989) *Proc. Natl Acad. Sci. USA* 86, 6773-6777.
- 52 Schneider, R. and Schweiger, M. (1991) *Oncogene* 6, 1807-1811.
- 53 Wenger, R.H. et al. (1988) *Biochem. Biophys. Res. Commun.* 156, 389-395.
- 54 Braun, T. et al. (1991) *EMBO J.* 10, 1885-1890.
- 55 Suzuki, N. et al. (1990) *Proc. Natl Acad. Sci. USA* 87, 8711-8715.
- 56 Krantz, D.D. et al. (1991) *J. Biol. Chem.* 266, 16801-16807.
- 57 Perin, J.P. et al. (1987) *J. Biol. Chem.* 262, 13269-13272.
- 58 Aruffo, A. et al. (1990) *Cell* 61, 1303-1313.
- 59 Miyake, K. et al. (1990) *J. Exp. Med.* 172, 69-75.
- 60 Yamamoto, T. et al. (1984) *Cell* 39, 27-38.
- 61 Südhof, T.C. et al. (1985) *Science* 228, 815-822.
- 62 Esser, V. et al. (1988) *J. Biol. Chem.* 263, 13282-13290.
- 63 Shevach, E.M. and Korty, P.E. (1989) *Immunol. Today* 10, 195-200.
- 64 Williams, A.F. (1991) *Cell Biol. Int. Reports* 15, 769-777.
- 65 LeClair, K.P. et al. (1986) *EMBO J.* 5, 3227-3234.
- 66 Deckert, M. et al. (1992) *J. Immunol.* 148, 672-677.
- 67 Williams, A.F. et al. (1988) *Immunogenetics* 27, 265-272.
- 68 Orr, H.T. et al. (1979) *Biochemistry* 18, 5711-5720.
- 69 Brown, J.H. et al. (1988) *Nature (London)* 332, 845-850.
- 70 Simister, N.E. and Mostov, K.E. (1989) *Nature (London)* 337, 184-187.
- 71 Beck, S. and Barrell, B.G. (1988) *Nature (London)* 331, 269-272.
- 72 Bjorkman, P.J. and Parham, P. (1990) *Annu. Rev. Biochem.* 59, 253-288.
- 73 Lawlor, D.A. et al. (1990) *Annu. Rev. Immunol.* 8, 23-63.
- 74 Leung, J.O. et al. (1985) *J. Biol. Chem.* 260, 12523-12527.
- 75 Sehgal, A. et al. (1988) *Mol. Cell. Biol.* 8, 3160-3167.

- 76 Strosberg, A.D. (1991) *Eur. J. Biochem.* 194, 1-10.
- 77 Dohlman, H.G. et al. (1991) *Adv. Biol. Sci.* 100, 1-10.
- 78 Freeman, M. et al. (1990) *Proc. Natl Acad. Sci. USA* 87, 1000-1004.
- 79 Freeman, M. et al. (1991) *Proc. Natl Acad. Sci. USA* 88, 1000-1004.
- 80 Huang, H.J. et al. (1987) *Proc. Natl Acad. Sci. USA* 84, 1000-1004.
- 81 Reth, M. (1989) *Nature (London)* 337, 184-187.
- 82 Letourneur, F. and Klausner, R.D. (1990) *Cell* 61, 1303-1313.
- 83 Irving, B.A. and Weiss, A. (1990) *Cell* 61, 1303-1313.
- 84 Rebbe, N.F. et al. (1991) *Proc. Natl Acad. Sci. USA* 88, 1000-1004.
- 85 Grundmann, U. et al. (1990) *Proc. Natl Acad. Sci. USA* 87, 1000-1004.
- 86 Andria, M.L. et al. (1991) *J. Immunol.* 146, 1000-1004.
- 87 Wright, M. et al. (1993) *Manuscript*.
- 88 Levy, S. et al. (1991) *J. Biol. Chem.* 266, 1000-1004.
- 89 Oren, R. et al. (1990) *Mol. Cell. Biol.* 10, 1000-1004.
- 90 Hupp, K. et al. (1989) *J. Immunol.* 143, 1000-1004.
- 91 Rubien, J.K. et al. (1989) *In L. University Press, Oxford*, pp. 1000-1004.
- 92 Rudd, C.E. et al. (1989) *Immunol. Rev.* 100, 1000-1004.
- 93 Samelson, L.E. et al. (1990) *Proc. Natl Acad. Sci. USA* 87, 1000-1004.
- 94 Reth, M. et al. (1991) *Immunol. Rev.* 100, 1000-1004.
- 95 Burkhardt, A.L. et al. (1991) *Proc. Natl Acad. Sci. USA* 88, 1000-1004.
- 96 Ullrich, A. and Schlessinger, J. (1990) *Cell* 61, 1303-1313.
- 97 Cantley, L.C. et al. (1991) *Cell* 61, 1303-1313.
- 98 Hanks, S.K. et al. (1988) *Science* 241, 1000-1004.
- 99 Yarden, Y. and Ullrich, A. (1990) *Cell* 61, 1303-1313.
- 100 Thomas, M.L. et al. (1985) *Cell* 41, 1000-1004.
- 101 Charbonneau, H. et al. (1988) *Biochemistry* 27, 1000-1004.
- 102 Tonks, N.K. et al. (1988) *Biochemistry* 27, 1000-1004.
- 103 Streuli, M. et al. (1990) *EMBO J.* 9, 1000-1004.
- 104 Streuli, M. et al. (1988) *J. Exp. Med.* 168, 1000-1004.
- 105 Streuli, M. et al. (1989) *Proc. Natl Acad. Sci. USA* 86, 1000-1004.
- 106 Hall, L.R. et al. (1988) *J. Immunol.* 140, 1000-1004.
- 107 George, D.G. et al. (1990) *Mol. Cell. Biol.* 10, 1000-1004.
- 108 Williams, A.F. and Barclay, A.P. (1988) *Immunogenetics* 27, 1000-1004.
- 109 Alt, F.W. and Rabbitts, T.H. (1984) *Cell* 39, 1000-1004.
- 110 Holmes, W.E. et al. (1991) *Science* 254, 1000-1004.
- 111 Gaugitsch, H.W. et al. (1990) *Cell* 61, 1000-1004.
- 112 Rippmann, F. et al. (1991) *Eur. J. Biochem.* 194, 1000-1004.

- 76 Strosberg, A.D. (1991) *Eur. J. Biochem.* 196, 1-10.
- 77 Dohlman, H.G. et al. (1991) *Annu. Rev. Biochem.* 60, 653-688.
- 78 Freeman, M. et al. (1990) *Proc. Natl Acad. Sci. USA* 87, 8810-8814.
- 79 Freeman, M. et al. (1991) *Proc. Natl Acad. Sci. USA* 88, 4931-4935.
- 80 Huang, H.J. et al. (1987) *Proc. Natl Acad. Sci. USA* 84, 204-208.
- 81 Reth, M. (1989) *Nature (London)* 338, 383-384.
- 82 Letourneur, F. and Klausner, R.D. (1992) *Science* 255, 79-82.
- 83 Irving, B.A. and Weiss, A. (1991) *Cell* 64, 891-901.
- 84 Rebbe, N.F. et al. (1991) *Proc. Natl Acad. Sci. USA* 88, 5192-5196.
- 85 Grundmann, U. et al. (1990) *DNA Cell Biol* 9, 243-250.
- 86 Andria, M.L. et al. (1991) *J. Immunol.* 147, 1030-1036.
- 87 Wright, M. et al. (1993) Manuscript in preparation.
- 88 Levy, S. et al. (1991) *J. Biol. Chem.* 266, 14597-14602.
- 89 Oren, R. et al. (1990) *Mol. Cell. Biol.* 10, 4007-4015.
- 90 Hupp, K. et al. (1989) *J. Immunol.* 143, 3787-3791.
- 91 Rubien, J.K. et al. (1989) In *Leukocyte Typing IV*, (Knapp, W. et al., eds) Oxford University Press, Oxford, pp.51-54.
- 92 Rudd, C.E. et al. (1989) *Immunol. Rev.* 111, 225-266.
- 93 Samelson, L.E. et al. (1990) *Proc. Natl Acad. Sci. USA* 87, 4358-4362.
- 94 Reth, M. et al. (1991) *Immunol. Today* 12, 196-201.
- 95 Burkhardt, A.L. et al. (1991) *Proc. Natl Acad. Sci. USA* 88, 7410-7414.
- 96 Ullrich, A. and Schlessinger, J. (1990) *Cell* 61, 203-212.
- 97 Cantley, L.C. et al. (1991) *Cell* 64, 281-302.
- 98 Hanks, S.K. et al. (1988) *Science* 241, 42-52.
- 99 Yarden, Y. and Ullrich, A. (1988) *Annu. Rev. Biochem.* 57, 443-478.
- 100 Thomas, M.L. et al. (1985) *Cell* 41, 83-93.
- 101 Charbonneau, H. et al. (1988) *Proc. Natl Acad. Sci. USA* 85, 7182-7186.
- 102 Tonks, N.K. et al. (1988) *Biochemistry* 27, 8695-8701.
- 103 Streuli, M. et al. (1990) *EMBO J.* 9, 2399-2407.
- 104 Streuli, M. et al. (1988) *J. Exp. Med.* 168, 1523-1530.
- 105 Streuli, M. et al. (1989) *Proc. Natl Acad. Sci. USA* 86, 8698-8702.
- 106 Hall, L.R. et al. (1988) *J. Immunol.* 141, 2781-2787.
- 107 George, D.G. et al. (1990) *Methods Enzymol.* 183, 333-351.
- 108 Williams, A.F. and Barclay, A.N. (1989) In *Immunoglobulin Genes*, (Honjo, T., Alt, F.W. and Rabbitts, T.H. eds) Academic Press, London, pp.361-387.
- 109 Holmes, W.E. et al. (1991) *Science* 253, 1278-1280.
- 110 Gaugitsch, H.W. et al. (1991) *Eur. J. Immunol.* 21, 377-383.
- 111 Flajnik, M.F. et al. (1991) *Immunogenetics* 33, 295-300.
- 112 Rippmann F. et al. (1991) *EMBO J* 10, 1053-1059.